# Summary Notes
# EE425X - Machine Learning: A Signal Processing Perspective

### Namrata Vaswani
### Iowa State University

These notes are a work in progress.For the latest version at any time, see the link `https://www.dropbox.com/scl/fi/nd7buqtksc55lei8e1k6z/ML_algorithms.pdf?rlkey=8e3oczn35ybwulvb82ecu5fcm&dl=0`

## Acknowledgement

## Contents

# Notation

In these notes $'$ and $^T$ and $^\top$ are all used to denote vector or matrix transpose. At a few other places, MATLAB notation is used too.

# 1 Math basics

- Probability - See file 322-recap- and cs229-

- Linear Algebra - See cs229-linalg

# 2 Monte Carlo methods and Data simulation details

## 2.1 Basics

- i.i.d. means independent and identically distributed. $X_1, X_2, \ldots, X_n$ are i.i.d. means

$$F_{X_1, X_2, \ldots, X_n}(x_1, x_2, \ldots, x_n) = \prod_{i=1}^{n} F_{X_i}(x_i) = \prod_{i=1}^{n} F_{X_1}(x_i)$$

and the same holds for their PDF/PMF also.

- $var(Y) = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$. For a sum of $n$ i.i.d. r.v.s $X_i$, $var(\sum_i X_i/n) = var(X_1)/n$

- For a unif(0,1), r.v., the CDF
$$F_X(x) = x, \ if \ 0 \le x < 1,$$

## 2.2 Monte Carlo for computing expected values

Let $Z$ be a vector or scalar or any set of r.v.s and let $g(Z)$ be a scalar function of $Z$.

In most cases
$$Eg := \mathbb{E}[g(Z)] = \int_z g(z) f_Z(z) dz$$

cannot be computed in closed form. One way to compute it is by Monte Carlo method. The idea is the following

- Generate $Z_i \sim f_Z(z)$, $i = 1, 2, ..N$ (generate $Z_i$ i.i.d. according to $f_Z(z)$)

- Compute the M-C estimate
$$\bar{g} := \frac{1}{N} \sum_i g(Z_i)$$

Can show that $\bar{g} \approx Eg$ with high probability if $N$ is large enough. This is due to law of large numbers.

Let
$$\sigma_g^2 = var(g(Z)) = \mathbb{E}[g^2(Z)] - (Eg)^2$$

For the M-C estimate,
$$\mathbb{E}[\bar{g}] = Eg, \quad \text{and} \quad var(\bar{g}) = \frac{1}{N}\sigma_g^2$$

where the variance expression uses the fact that the $Z_i$s are independent

Chebyshev's inequality: For a r.v. $Z$ with mean $\mu$, and finite variance $\sigma^2$, we have

$$\Pr(|Z - \mu| > \epsilon) \le \frac{\sigma^2}{\epsilon^2}$$

This is one of the simplest and most applicable law of large numbers inequalities.

Using Chebyshev inequality, and the above expressions for mean and variance of $\bar{g}$ we can show that

$$\Pr(|\bar{g} - Eg| > \epsilon) \leq \frac{\sigma_g^2}{N\epsilon^2}$$

Typically we pick $\epsilon$ to be a small fraction of the quantity being estimated, e.g., we can let $\epsilon = 0.1Eg$. Then,

$$\Pr(|\bar{g} - Eg| > 0.1Eg) \leq \frac{100}{N}\frac{\sigma_g^2}{(Eg)^2}$$

Thus if

$$N \geq \frac{100}{p_0}\frac{\sigma_g^2}{(Eg)^2}$$

then, with probability at least $1 - p_0$, $|\bar{g} - Eg| < 0.1Eg$.

Qualitative conclusion: the number of samples required to get accurate estimates with high probability (w.h.p.), is proportional to $\frac{\sigma_g^2}{(Eg)^2}$. The larger the variability in $g(Z)$, the more the samples we need.

We can approximate $\sigma_g^2$ also using M-C, but that approx uses a choice of $N$. A hit and trial procedure can be tried and works sometime.

## 2.3   Monte Carlo in ML

We use M-C to calculate the generalization error in ML. Consider a ML problem. Suppose we have assumed a model and learned its parameters. Our goal is to decide how good our learned model is at predicting unseen data. We split available data into a training set and a test set. Learn the model parameters using the training set. Compute test/generalization error on the test set.

For a loss function $\ell(\boldsymbol{x}, y)$, we want to compute

$$GeneralizationError = \mathbb{E}[\ell(\boldsymbol{x}, y)]$$

for test data $\boldsymbol{x}$ and output $y$ generated by nature. We approximate this using M-C as

$$\frac{1}{m_{test}}\sum_i \ell(\boldsymbol{x}_i^{test}, y_i^{test})$$

A common loss function example is squared prediction error $\ell(\boldsymbol{x}, y) = (y - h_{\hat{\theta}}(\boldsymbol{x}))^2$; here $h_\theta(\boldsymbol{x})$ is the assumed model, we learn its parameter vector $\theta$ using training data and denote the learned estimate by $\hat{\theta}$

Observe: in the above we do not need to know the distribution of $\boldsymbol{x}$.

## 2.4   Random variable generation

MATLAB and NumPy within Python have the functions

- rand: generate a r.v. that is uniformly distributed between 0 and 1

- randn: generate a r.v. that is standard Gaussian (Gaussian with mean 0 and variance 1)

- rand(n,1): generate $n$ i.i.d. unif(0,1) r.v.s as a vector

- randn(n,1): generate $n$ i.i.d. N(0,1) r.v.s as a vector, or said another way $\boldsymbol{x} = randn(n, 1) \sim N(0, \boldsymbol{I})$.

### 2.4.1 How to generate a discrete r.v. $Y$ with a finite number of possible values from a unif(0,1) r.v.

Suppose $Y$ takes three possible values, say 1, 2, 3, and $p_Y(1) = p_1, p_Y(2) = p_2, p_Y(3) = p_3 = 1 - p_1 - p_2$. One way to generate $Y$ is as follows:

- Generate $X = rand$

- If

  - $X < p_1$, then $Y = 1$, else-if
  - $p_1 \leq X < p_1 + p_2$, then $Y = 2$ else
  - $Y = 3$

Proof of why above works:

$$\Pr(Y = 1) = Pr(X < p_1) = p_1, \ \Pr(Y = 2) = Pr(p_1 \leq X < p_1 + p_2) = (p_1 + p_2) - p_1 = p_1$$

### 2.4.2 How to generate a new continuous r.v. $Y$ with CDF $F_Y(y)$ from a unif(0,1) r.v.

- Generate $X = rand$ is a unif(0,1) r.v.

- Compute $Y = F_Y^{-1}(X)$.

- Here $F_Y(y)$ be the CDF of $Y$. Since $Y$ is a continuous r.v. this is a strictly increasing function, and so it is invertible.

Proof of why above works:
$$\Pr(F_Y^{-1}(X) \leq z) = \Pr(X \leq F_Y(z)) = F_Y(z)$$

Last inequality uses CDF of unif(0,1). Thus, $Y$ indeed has CDF $F_Y$.

### 2.4.3 How to generate a $N(\mu, \Sigma)$ r vector

- Generate $\boldsymbol{z} = randn(n, 1)$

- Compute matrix square root of $\Sigma$ using Cholesky or SVD: find a $B$ so that $BB^\top = \Sigma$

- Generate $\boldsymbol{y} = B\boldsymbol{z} + \mu$.

## 2.5 Data Simulation

- Why simulate: It allows us to check why the prediction error on real data is large. Is it because our assumed model is a bad approximation of real data, or is it because the learning algorithm for our assumed model is bad? Simulation helps fix the second problem. The learning algorithm may be not good (if you are designing a new one), or you may have coding errors due to silly mistakes. Data simulation helps check that.

How to simulate:

- Assume a model on how $\boldsymbol{x}$ relates to $y$, and include a modeling error term (usually additive):

$$y = h_\theta(\boldsymbol{x}) + e$$

- Assume a value for the parameter vector $\theta$.

- Assume a PDF/PMF of the modeling error or noise $e$.

- Assume a PDF/PMF on $\boldsymbol{x}$

- Create training dataset: For $i = 1, 2, \ldots, m$, Generate

    - $\boldsymbol{x}_i$ i.i.d. from the assumed PDF/PMF
    - $e_i$ i.i.d. from the assumed PDF/PMF
    - Create $y_i = h_\theta(\boldsymbol{x}_i) + e_i$

- Create test dataset: proceed exactly as above.

Learn the parameters (estimation)

- Find $\hat{\theta} \in \arg\min_{\tilde{\theta}} \sum_i l(\boldsymbol{y}_i, h_{\tilde{\theta}}(\boldsymbol{x}_i))$

Compute the errors: always compute normalized errors

- learning error in learning $\theta$ (only possible for simulated data where true $\theta$ is known):

$$||\hat{\theta} - \theta||_2^2 / ||\theta||_2^2$$

- generalization error / prediction error on test data

$$\sum_i \ell(\boldsymbol{y}_i^{test}, h_{\hat{\theta}}(\boldsymbol{x}_i^{test})) / \text{normalize-appropriately}$$

# 3 Data Simulation – more

## 3.1 Why simulate – why test code on simulated data first?

Consider the house price prediction based on house features example. In the previous section, we said we use linear regression to model the data and predict the price. If I write code to learn theta and apply it to a real dataset directly and suppose it does not "work too well" (my code does not give very good predictions on test data). How do I know if (i) the linear regression model is wrong or (ii) my learning approach (normal equations or GD) is wrong, or (iii) there is a code bug (I have an extra factor of 2 at some place by mistake)?

A partial fix to the above problem to above it is to first simulate data using the linear regression model. So when I test my learning code on this data I know what is true $\theta$ I am looking for. Then I can try to fix (ii) and (iii) issues. Once these are fixed, then we try the same code on real data and then maybe compare with another model to check which one is better.

## 3.2 Monte Carlo Estimation of Expected values

Suppose the goal is to compute/approximate $\mathbb{E}[z]$ (expected value or average). Suppose also we know the probability distribution of $z$, $p(z)$, and we have a way to generate (simulate) independent realizations of $z$. Then we can approximate

$$\mathbb{E}[z] \approx \frac{1}{m} \sum_{i=1}^{m} z^{(i)}, \ z^{(i)} \stackrel{\text{i.i.d.}}{\sim} p(z)$$

We use this to approximate prediction error in ML among other things.

## 3.3 Understanding a multivariate Gaussian distribution

See these two files `http://cs229.stanford.edu/section/gaussians.pdf`
`http://cs229.stanford.edu/notes2020fall/notes2020fall/cs229-prob.pdf`

### 3.4 How to simulate

Task: Generate your own data to simulate the linear regression model $y = \boldsymbol{\theta}^T \boldsymbol{x} + e$.

Recall: $n = d + 1$, thus $\theta$ is a vector of length $n$. The $\boldsymbol{x}^{(i)}$s are of length $d$ with 1 appended as first entry.

As an example, suppose $d = 5$, then $n = 6$. Can set $\theta = [10, 1, -1, -3, 4, 2]^\top$.

Generate $m$ such independent training data vectors. Also generate another set of $m_{test}$ independent data vectors for the testing step.

Let $\mathcal{N}(\mu, \boldsymbol{\Sigma})$ denote the Gaussian distribution with mean $\mu$ and covariance $\boldsymbol{\Sigma}$.

Do this as follows:

(1) Fix $\boldsymbol{\theta}$ once as an $n$ length vector.

(2) Training data generate:

For $i = 1$ to $m$, generate $\boldsymbol{x}^{(i)} \sim \mathcal{N}(0, I_d)$, $e^{(i)} \sim \mathcal{N}(0, \sigma_e^2)$ and $y^{(i)} = \boldsymbol{\theta}^T [1, \boldsymbol{x}^{(i)}] + e^{(i)}$.

(3) Test data generate:

In a different Loop, repeat (2) for $i = 1$ to $m_{test}$. Use this data in the testing step and not for training.

No-for-loop version of (2):

$\boldsymbol{X} = [ones(m, 1), randn(m, d)]$ , $\boldsymbol{e} = \sigma_e * randn(m, 1)$, $\boldsymbol{y} = \boldsymbol{X}\theta + \boldsymbol{e}$.

#### 3.4.1 Notes

*For generating the data, I have suggested using the Gaussian distribution just as an example. But you do not have to use the Gaussian. You could also use any other distribution, e.g., the uniform distribution.*

On the other hand, for the error $e^{(i)}$, we have assumed that it is Gaussian in our model (that is why the squared loss is justified). So for generating error, you *do have to* use the Gaussian distribution.

### 3.5 Estimation error computations: how to know your code is correct

- Always compute normalized errors: if code works well error should be less than 1. When increase $m$ (number of samples), while keeping everything else fixed, this error should reduce.

- If data were simulated, then one can compute (normalized) parameter estimation error $\frac{\|\theta - \hat{\theta}\|_2^2}{\|\theta\|_2^2}$

- In all cases, one can compute (normalized) prediction error / training loss: $\frac{(1/m) \sum_{i=1}^{m} (\boldsymbol{y}^{(i)} - \hat{\boldsymbol{y}}^{(i)})^2}{(1/m) \sum_{i=1}^{m} (\boldsymbol{y}^{(i)})^2}$

- Normalized Test MSE / Generalization Error computation:

  - Simulated data: Generate an independent set of test data
  - Real data: Two options
    * Hold-out Cross Validation: Split available data into disjoint training and test sets, use most data for training, rest for testing.
    * Or, use Leave-One-Out Cross-Validation, see Sec. 4
  - For each test data point, compute $\hat{\boldsymbol{y}}_{test} = h_{\hat{\theta}}(\boldsymbol{x}_{test})$
  - Compute normalized average squared error $\frac{\sum_{i=1}^{m_{test}} (\boldsymbol{y}_{test}^{(i)} - \hat{\boldsymbol{y}}_{test}^{(i)})^2}{\sum_{i=1}^{m_{test}} (\boldsymbol{y}_{test}^{(i)})^2}$

- For different problems, squared error can be replaced by other error measures.

## 4 Cross Validation for Supervised Learning: training and test datasets

In case of supervised learning, we need to define training and test datasets. These could be non-overlapping (easiest to code, analyze, fastest to run) or overlapping (needed if less total data).

## 4.1 Simple Commonly Used Approach: Hold-Out Cross Validation

Simplest way to define training and test datasets is to use a certain fraction of the data for training and the rest for testing. A common rule of thumb is 80-20 or 70-30. In this case training and test data are disjoint which makes it easier to impose statistical independence assumptions. Also this is easiest to code in and the code runs fastest.

## 4.2 Leave-one-out cross validation: any general task**

Use of distinct subsets is easiest to code in but wastes a lot of data. The other extreme solution is to use leave-one-out cross-validation. This means we loop over all the data multiple times. In the $i$-th loop, we use all but the $i$-th data point for training and use the $i$-th one for testing. Compute error on this $i$-th one and store it. At the end of this loop, average the stored errors from each loop to compute the average generalization error / predictor error.

- for $i = 1$ to $m_{tot}$

  - at iteration $i$, define the training data matrix $\boldsymbol{X}$ and vector $\boldsymbol{y}$ as follows: use *all* data except the $i$-th

$$\boldsymbol{X} = \begin{bmatrix} -(\boldsymbol{x}^{(1)})^T- \\ -(\boldsymbol{x}^{(2)})^T- \\ \vdots \\ -(\boldsymbol{x}^{(i-1)})^T- \\ -(\boldsymbol{x}^{(i+1)})^T- \\ \vdots \\ -(\boldsymbol{x}^{(m)})^T- \end{bmatrix}$$

$$\boldsymbol{y} = \begin{bmatrix} \boldsymbol{y}^{(1)} \\ \boldsymbol{y}^{(2)} \\ \vdots \\ \boldsymbol{y}^{(i-1)} \\ \boldsymbol{y}^{(i+1)} \\ \vdots \\ \boldsymbol{y}^{(m)} \end{bmatrix}$$

  - Learn the model parameters use $\boldsymbol{X}, \boldsymbol{y}$ as the training data
  - Use the $i$-th data points for computing the error $err(i)$.

  end for

Compute

$$TestMSE = \frac{1}{m_{tot}} \sum_{i=1}^{m_{tot}} err(i), NormalizedTestMSE = \text{as needed for the problem}$$

### 4.2.1 Leave-one-out cross validation for Linear regression

Given $m$ training data points $\boldsymbol{x}^{(i)}, \boldsymbol{y}^i i$, $i = 1, 2, \ldots, m$. Suppose the goal is to evaluate the validity of a linear regression model on this data. We will compute Normalized-Test-MSE as follows.

- for $i = 1$ to $m_{tot}$

– at iteration $i$, define the training data matrix $\boldsymbol{X}$ and vector $\boldsymbol{y}$ as follows: use *all* data except the $i$-th. This is defined above.

$$
\boldsymbol{X} = \begin{bmatrix} -(\boldsymbol{x}^{(1)})^T- \\ -(\boldsymbol{x}^{(2)})^T- \\ \vdots \\ -(\boldsymbol{x}^{(i-1)})^T- \\ -(\boldsymbol{x}^{(i+1)})^T- \\ \vdots \\ -(\boldsymbol{x}^{(m)})^T- \end{bmatrix}, \quad \boldsymbol{y} = \begin{bmatrix} \boldsymbol{y}^{(1)} \\ \boldsymbol{y}^{(2)} \\ \vdots \\ \boldsymbol{y}^{(i-1)} \\ \boldsymbol{y}^{(i+1)} \\ \vdots \\ \boldsymbol{y}^{(m)} \end{bmatrix}
$$

– Compute $\hat{\theta} = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{y}$
– $\hat{\boldsymbol{y}}^{(i)} = \hat{\theta}^T\boldsymbol{x}^{(i)}$
– Compute $err(i) = (\boldsymbol{y}^{(i)} - \hat{\boldsymbol{y}}^{(i)})^2$

end for

Compute

$$
TestMSE = \frac{1}{m_{tot}} \sum_{i=1}^{m} err(i), \; NormalizedTestMSE = TestMSE/(\sum_{i=1}^{m_{tot}} (\boldsymbol{y}^{(i)})^2)
$$

# 5 Optimization Algorithms (GD, SGD, and AltMin)

## 5.1 Gradient Descent (GD)

The GD approach is an iterative algorithm to find a local minimizer of a cost function. Which minimizer is found depends on how one initializes. It does not always converge and of course local minimizer may not be a global minimizer either. But if cost function is convex, it will converge to a global minimizer. Moreover if minimizer is also unique (cost is strictly convex), then the only correct solution can be found. An example of this is the squared loss for linear regression. It is convex always. It is strictly convex if $m \geq n$ and $\boldsymbol{X}$ has full rank $n$.

The GD Algorithm to minimize any cost $J(\theta)$ is as follows. Recall

$$
J(\theta) := \sum_{i=1}^{m} cost(\theta, \boldsymbol{x}^{(i)}, y^{(i)})
$$

where *cost* is squared error (as in previous section) or can be something else (as in later sections on Logistic Regression).

1. Initialize $\hat{\theta}$ as the zero vector
   (or anything else: this choice is problem dependent, for linear regression zero vector is good)

2. For $t = 1$ to $Tmax$ do
   (use a large value for $Tmax$)

   (a) $\hat{\theta}_{old} \leftarrow \hat{\theta}$
   (b) $\hat{\theta} \leftarrow \hat{\theta} - \mu \nabla_\theta J(\hat{\theta})$
   (c) Exit loop if $\|\hat{\theta} - \hat{\theta}_{old}\|_2/\|\hat{\theta}_{old}\|_2 \leq \epsilon$ with $\epsilon$ a very small tolerance, e.g., set $\epsilon = 10^{-6}$ (this value actually also depends on noise level)

   End for

## 5.2 Setting the step size $\mu$

- When using $J(\theta)$ as above (it is the sum of $m$ terms), then use $\mu = c/m$ with $c$ being a number between 0 and 1.

- Rule of thumb: Reduce $c$ it if the cost seems to not decrease at all or starts increasing. Increase it if the cost decreases but very slowly with iteration.

  - Detailed: start with $c = 0.5$, check if error goes down with iterations, if not it is too large. Use a smaller $c$ say $c = 0.05$. Check again. If still does not go down, use an even smaller step size.

- As you reduce $c$, you will need MORE total number of iterations.

- Good coding practice: Always use a $T_{max}$: maximum number of iterations and an exit criterion.

- Extra info: *Variable step size: Can also vary $\mu$ with iterations, start with a larger value but reduce with iterations.*

## 5.3 Stochastic GD (SGD) and Stochastic Mini-Batch GD **

$J(\theta)$ is typically an average of $m$ terms; in fact it always is under our assumption of different training data points being i.i.d. As a result, its gradient is also a sum of $m$ terms divided by $m$.

If $m$ is large, computing the full gradient at each iteration can be expensive. Also, sometimes not all data is available immediately.

Stochastic GD idea: sum over a subset of the $m$ gradients at each iteration. Pick this subset randomly or use other strategies.

Details: see videos mentioned next.

### 5.3.1 Main idea

Let

$$J(\theta) := \sum_{i=1}^{m} J_i(\theta; \boldsymbol{x}^{(i)}, y^{(i)})$$

- Initialize $\hat{\theta}$.

- Repeat

  - Randomly shuffle samples in the training set.
  - Loop through all points in the training set one at a time.
    * For $i = 1, 2, \ldots, m$, do:
    $$\hat{\theta} \leftarrow \hat{\theta} - \eta \nabla J_i(\hat{\theta})$$

  until stopping criterion reached

### 5.3.2 Good videos on stochastic and mini-batch SGD

Basic GD video `https://www.youtube.com/watch?v=yFPLyDwVifc`
    Stochastic mini-batch GD `https://www.youtube.com/watch?v=4qJaSmvhxi8`

## 5.4 Alternating Minimization (AltMin) algorithm

Suppose we need to solve

$$\min_{\theta} J(\theta)$$

AltMin is another approach to solve an optimization problem. It is a better one to use that Gradient Descent when the variable to be minimized over, $\theta$, can be split into two subsets of variables $\theta = [\theta_1, \theta_2]$ such that minimizing over $\theta_1$ keeping $\theta_2$ fixed and vice versa is either closed form or otherwise easy to do.

AltMin proceeds as follows

1. Randomly initialize $\theta_1$ to $\hat{\theta}_1$.

2. Repeat

    (a) Minimize over $\theta_2$ keeping $\theta_1$ fixed at $\hat{\theta}_1$, i.e., compute

    $$\hat{\theta}_2 = \arg\min_{\theta_2} J(\hat{\theta}_1, \theta_2)$$

    (b) Minimize over $\theta_1$ keeping $\theta_2$ fixed at $\hat{\theta}_2$, i.e., compute

    $$\hat{\theta}_1 = \arg\min_{\theta_1} J(\theta_1, \hat{\theta}_2)$$

    until "convergence" i.e. estimates of $\hat{\theta}_1$ do not change much from previous to current iteration.

Like GD, AltMin also only converges to the local minimum of the cost function. Thus to make it work better one can run it $N$ times with different random initializations. Pick the solution that result in the smallest cost function value.

# 6 Least Squares (LS) Estimation and Regularized LS

The goal is to recover $\theta$ from

$$\boldsymbol{y} := \boldsymbol{X}\theta + \boldsymbol{e}$$

where $\boldsymbol{e}$ is the modeling error or noise. The goal is to learn $\theta$ so that $||\boldsymbol{e}||^2$ is minimized.

## 6.1 Ordinary LS

The opt problem that needs to be solved to learn the parameter vector $\theta$ for Lin Reg is called the "Least Squares" problem. This involves solving

$$\min_{\theta} \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2$$

The above is called *Ordinary Least Squares* or *Ordinary Linear Regression*. This problem also occurs in various other domains. Many examples in communications and in signal processing. Thus it is useful to understand its various modifications.

## 6.2 Regularized LS – sparse features

ell-1 regularized LS: this ensures the vector $\theta$ is sparse (has many zero entries), larger $\lambda$ makes it sparser. Use if only a subset of all the features in the feature vector matter but do not know which ones.

$$\min_{\theta} \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2 + \lambda\|\theta\|_1$$

The above is also called *Lasso* or *ell-1 regularization*.
   **Do not try to write your own GD code for this. Use cvxopt or cvxpy etc.**
**Vary $\lambda$ as multiples of $m$, so $\lambda = \beta m$, with $\beta = 0.01, 0.1, 1, 10$ etc**

## 6.3 Regularized LS – $\ell_2$ norm regularization

We may claim that no house feature is way too important, and thus, no entry of $\theta$ should be too large. In this case we solve

$$\min_\theta \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2 + \lambda\|\theta\|^2$$

The above is also called *Ridge Regression* or *ell-2 regularization*

## 6.4 Regularization: Use of prior knowledge and Recursive LS

In many settings, we may have prior knowledge that $\theta$ should be close to $\theta_0$. For example, we may have $\theta_0$ from houses sold a year ago and their features. In such cases, we solve

$$\min_\theta \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2 + \lambda\|\theta - \theta_0\|^2$$

Recursive LS: consider house price example again and suppose we want to update our model each week. Instead of starting from scratch, can we just "update" the model: save time/computation?

Turns out we can using recursive LS: suppose we update with each new house sold, then we get this

$$\hat{\theta}_0 = \theta_0, P_0 = \lambda\boldsymbol{I},$$

$$K_t = P_{t-1}\boldsymbol{x}^{(i)}(\lambda + \boldsymbol{x}^{(i)\top}P_{t-1}\boldsymbol{x}^{(i)})^{-1},$$

$$P_t = (\boldsymbol{I} - K_t\boldsymbol{x}^{(i)\top})P_{t-1},$$

$$\hat{\theta}_t = \hat{\theta}_{t-1} + K_t(\boldsymbol{y}^{(i)} - \boldsymbol{x}^{(i)\top}\hat{\theta}_{t-1}),$$

Suppose we update after a set of $m_t$ new houses sold (so $\boldsymbol{x}^{(i)}$ is an $n \times m_t$ matrix and $\boldsymbol{y}^{(i)}$ is a $m_t \times 1$ vector), then we get this

$$\hat{\theta}_0 = \theta_0, P_0 = \boldsymbol{I},$$

$$K_t = P_{t-1}\boldsymbol{x}^{(i)}(\lambda\boldsymbol{I} + \boldsymbol{x}^{(i)\top}P_{t-1}\boldsymbol{x}^{(i)})^{-1},$$

$$P_t = (\boldsymbol{I} - K_t\boldsymbol{x}^{(i)\top})P_{t-1},$$

$$\hat{\theta}_t = \hat{\theta}_{t-1} + K_t(\boldsymbol{y}^{(i)} - \boldsymbol{x}^{(i)\top}\hat{\theta}_{t-1}),$$

# 7 Supervised Learning, Unsupervised Learning, Model-based Learning

Given observed data (or features of the observed data) or other "input" $\boldsymbol{x}$, the goal in ML is to predict/compute some function of $\boldsymbol{x}$ that is usually denoted by $y$ (often called "output"). In most of this course, $\boldsymbol{x}$ is an $n \times 1$ vector and $y$ is a scalar. For example, $\boldsymbol{x}$ can be the feature vector of key attributes of a house, while $y$ can be its price. In this case both are real valued. Or $y$ can be a binary decision about whether a buyer buys the house or not. As a different example, $\boldsymbol{x}$ can be a vectorized image while $y$ is the class label for the image (dog vs cat or dog vs cat vs human etc).

In supervised learning, we first decide a modeling strategy to model the input-output relationship; then come up with an algorithm to "learn" parameters given the training data (which is a set of $m$ input-output pairs). All of this is done so that the "learnt model" can be used to predict $y$ (get $\hat{y}$) for a new query $\boldsymbol{x}$. "Predict" is often also called "estimate" (if $y$ is real-valued) and it is also called "detect" or "classify" (if $y$ is binary/discrete-valued).

Learning algorithms can be supervised or unsupervised. In supervised learning, we are provided with "training data" that allows us to "learn" the parameters used by the model that our algorithm relies on. Goal is to predict $y$ using observed data or features $\boldsymbol{x}$.

$\boldsymbol{x}$ is $n \times 1$, $y$ is a scalar. We use $\theta$ to denote the set of parameters used by our assumed model.

In many settings, the assumed model that predicts $y$ is denoted $h_\theta(\boldsymbol{x})$. Since the model is never perfect, we assume that the "true" output $y$ satisfies

$$y = h_\theta(\boldsymbol{x}) + e$$

where $e$ is the *modeling error or noise.* This is typically modeled as a random variable with a probability density function (PDF), typically zero mean Gaussian and independent and identically distributed (i.i.d.) in each new sample.

Training data consists of $m$ input/output pairs $\{\boldsymbol{x}^{(i)}, y^{(i)}\}, i = 1, 2, \ldots, m$. The modeling error / noise $e$ We use these to "learn" $\theta$. Once that is done, we can predict $y$ from $\boldsymbol{x}$ using the above equation.

## 7.1 Examples

Examples of supervised learning include: linear regression, classification - logistic regression or GDA or naive-bayes/spam-filter, SVMs, deep learning.

In case of unsupervised learning, we are just provided data and we need to "make sense" of it. This can mean different things. In case of PCA this means finding a lower dimensional subspace to represent the data. In case of clustering, this means partition the data into $K$ disjoint "clusters"; $K$ itself may be known or unknown.

In Model-based Learning, the observed data model is provided (usually the Physics of the problem defines this). Examples include least squares (LS) problems, Bayesian LS, Bayesian estimation, or structured data recovery problems such as sparse recovery (compressive sensing) or low-rank matrix recovery. Standard LS is usually the oversampled data setting. Bayesian LS or structure models are used in case of undersampled data - some data is missing, corrupted, or deliberately undersampled.

The model learning (regression coefficients' vector learning) step of linear regression is an instance of LS. If limited training data is available, and we model the regression vector as being sparse, this can become an instance of sparse recovery.

Unsupervised deep learning, which is a newish area, in which one uses a single training data (that is a sequence of vectors or images or image Fourier transforms) to both learn the deep model parameters and predict the output - falls in the second or third category. Used in accelerated dynamc MRI applications.

## 7.2 Supervised Learning pipeline: probabilistic models

All learning assumes a model. Supervised learning also assumes availability of training data. Supervised learning itself can be probabilistic or not (SVM is not for example). Goal of learning: given input (feature vector) $\boldsymbol{x}$, predict output $y$.

- Decide Model:
  We first decide a modeling strategy to model the input-output relationship: how $\boldsymbol{x}$ and $y$ are related

  - in classical (non-Bayesian) learning, this is specified by $p(y; \boldsymbol{x}, \theta)$

  $$LinReg: \ p(y; \boldsymbol{x}, \theta) = \mathcal{N}(y; \theta^\top \boldsymbol{x}, \sigma_e^2)$$

  $$(Binary)LogReg: \ p(y; \boldsymbol{x}, \theta) = g(\theta^\top \boldsymbol{x})^y [1 - g(\theta^\top \boldsymbol{x})]^{1-y}$$

  $$(MultiClass)LogReg:$$

  - in generative (Bayesian) learning, this is specified by $p(y; \theta), p(\boldsymbol{x}|y; \theta)$ which can then specify $p(y|\boldsymbol{x})$ (Bayes rule)

  $$GDA: p(y; \theta) = \prod_{k=1}^{K-1} \phi_k^{[y==k]}, \ p(\boldsymbol{x}|y; \theta) = \mathcal{N}(\boldsymbol{x}; \mu_y, \Sigma)$$

- Learning:
  We use training data set $\{\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}\}_{i=1}^{m}$ to learn the model parameters $\theta$. Arrange these into matrix $\boldsymbol{X}$ and vector $\boldsymbol{y}$ Learning is done using Max Likelihood Estimation (MLE).

  – Classical case: MLE finds

  $$\arg\max_{\theta} p(\boldsymbol{y}; \boldsymbol{X}, \theta) := \prod_{i=1}^{m} p(\boldsymbol{y}^{(i)}; \boldsymbol{x}^{(i)}, \theta)$$

  – Bayesian case: MLE finds

  $$\arg\max_{\theta} p(\boldsymbol{y}, \boldsymbol{X}; \theta) := \prod_{i=1}^{m} p(\boldsymbol{y}^{(i)}, \boldsymbol{x}^{(i)}; \theta) = \prod_{i=1}^{m} p(\boldsymbol{y}^{(i)}; \theta) p(\boldsymbol{x}^{(i)} | \boldsymbol{y}^{(i)}; \theta)$$

- Prediction:
  Given a query $\boldsymbol{x}$, we use the learned model parameters, $\hat{\theta}$, and the model to predict $\hat{y}$

  – Classical case: Max Likelihood for finding $y$

  $$\hat{y} = \arg\max_{y} p(y; \boldsymbol{x}, \hat{\theta})$$

  – Bayesian case: Max A Posteriori (MAP) rule for finding $y$

  $$\hat{y} = \arg\max_{y} p(y | \boldsymbol{x}; \hat{\theta}) = \arg\max_{y} [p(y; \hat{\theta}) p(\boldsymbol{x} | y; \hat{\theta})]$$

- IMPORTANT: finding the argument to maximize a function is the same as finding the argument that minimizes its negative logarithm. Always use $-\log(.)$ for products of PMFs or PDFs - simpler expressions for deriving closed-form solutions; and this is also needed to prevent overflow/underflow in computer code.

  To add: PDF, PMF difference and similarity, Bayes rule; Examples of each kind - Lin Reg, Log Reg, GDA, Spam filter

## 7.3 Supervised Learning: not probabilistic models

- SVM

# 8 Supervised Learning: Linear Regression

In the setting we have talked about in class, $\boldsymbol{x}$ is a real-valued $n \times 1$ vector and $y$ is a real-valued scalar [1]. The parameter vector $\theta$ is also an $n \times 1$ vector

## 8.1 Model

In linear regression, $h_\theta(\boldsymbol{x})$ is a linear function of $\boldsymbol{x}$.

$$h_\theta(\boldsymbol{x}) = \theta^T \boldsymbol{x}.$$

(I sometimes may use $'$ for transpose – MATLAB notation)

---

[1] In more general settings $y$ can also be a real-valued vector (this will not be discussed in our class).

## 8.2  Including a nonzero mean in the model

To include a non-zero mean in the model, one can replace both $\boldsymbol{x}$ and $\theta$ by $n+1$ length vectors as follows. Let

$$\tilde{\theta} = \begin{bmatrix} \theta_0 \\ \theta \end{bmatrix}$$

and

$$\tilde{\boldsymbol{x}} = \begin{bmatrix} 1 \\ \boldsymbol{x} \end{bmatrix}$$

and we let

$$h_{\tilde{\theta}}(\tilde{\boldsymbol{x}}) = \tilde{\theta}^T \tilde{\boldsymbol{x}}.$$

With this model, we do everything explained above using $\tilde{\boldsymbol{x}}$ and $\tilde{\theta}$ to replace $\theta$ and $\boldsymbol{x}$ respectively.

<span style="color:red">Here and below, we let $\theta$ be an $n = d + 1$ length vector and same for $\boldsymbol{x}^{(i)}$ (append 1 as its first entry).</span>

## 8.3  Learning $\theta$: minimize squared loss

<span style="color:red">Here and below, we let $\theta$ be an $n = d + 1$ length vector and same for $\boldsymbol{x}^{(i)}$ (append 1 as its first entry).</span>

The most common approach to learn $\theta$ is to assume a squared error loss and try to minimize it, i.e., find

$$\arg\min_{\theta} J(\theta) := \sum_{i=1}^{m} (y^{(i)} - \theta^\top \boldsymbol{x}^{(i)})^2.$$

Define an $m \times n$ input data matrix $\boldsymbol{X}$ with $[1, (\boldsymbol{x}^{(i)})^T]$ as its rows, i.e., let

$$\boldsymbol{X} = \begin{bmatrix} 1, & -(\boldsymbol{x}^{(1)})^T - \\ 1, & (\boldsymbol{x}^{(2)})^T - \\ \vdots, & \vdots \\ 1, & -(\boldsymbol{x}^{(m)})^T - \end{bmatrix} \tag{1}$$

and define an $m \times 1$ vector $\boldsymbol{y}$ with $y^{(i)}$ as its columns.

Then, $J(\theta)$ can be expressed more compactly as

$$J(\theta) := \sum_{i=1}^{m} (y^{(i)} - \theta^T \boldsymbol{x}^{(i)})^2 = \|\boldsymbol{y} - \boldsymbol{X}\theta\|_2^2$$

## 8.4  Solutions for minimizing squared loss: also called Least Squares (LS) Estimation

1. We can get a closed form solution by taking the derivative of $J(\theta)$ w.r.t. $\theta$ and setting it to zero. When $\boldsymbol{X}$ is full rank $n$ (a necessary condition for this is $m \geq n$), this simplifies to

$$\hat{\theta} = (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \boldsymbol{y}$$

*Details:* $\nabla J(\theta) = 2\boldsymbol{X}^\top (\boldsymbol{y} - \boldsymbol{X}\theta)$. Set this equal to zero and solve for $\hat{\theta}$. Reason: the gradient is zero at every minimizer, maximizer or stationary point of a differentiable cost function. gradient = 0 is a necessary condition for a point to be a minimizer. in this case, since $J(.)$ is strongly convex it has exactly one stationary point which is THE minimizer. Thus setting $\nabla J(\theta) = 0$ helps us find THE minimizer.

2. For an $m \times n$ matrix with $m \geq n$, $(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top$ is the pseudo-inverse of $\boldsymbol{X}$, denoted $\boldsymbol{X}^\dagger$. Thus, we can also write the above solution as

$$\hat{\theta} = \boldsymbol{X}^\dagger \boldsymbol{y}$$

*Extra details:* The pseudo-inverse can be computed more efficiently than in the formula above: this is computed by first computing the singular value decomposition (SVD) of the matrix $\boldsymbol{X}$. Suppose the SVD of $\boldsymbol{X} = U\Sigma V^\top$ where $U$ is an $m \times m$ unitary matrix. $V$ is an $n \times n$ unitary matrix and $\Sigma$ is diagonal with non-negative entries. Then $\boldsymbol{X}^\dagger = V\Sigma^\dagger U^T$. For a rectangular diagonal matrix such as $\Sigma$, we get the pseudo-inverse by taking the reciprocal of each non-zero element on the diagonal, leaving the zeros in place, and then transposing the matrix.

*Do not use the pinv function of numpy though. That gives wrong output for $m = n$.*

If algorithms for computing matrix inverse or matrix pseudo-inverse were "exact" (and not iterative) then the above two approaches would return the exact same solution. But they are not. Thus, when $\boldsymbol{X}$ is "well-conditioned", both approaches above should return same solution, but not otherwise.

3. For large sized problems (where $n$ is large), using Gradient Descent (GD) is a better idea. Since problem is convex, GD should, in principle, converge to above solution starting from any initialization, and should converge pretty quickly. We explain this in Sec. 5.1.

4. *Approximate but even faster solution:* Stochastic GD (S-GD) can be used. Advantages: faster per iteration; needs lesser memory; and is useful to get a fully streaming algorithm. But no easy guarantees on whether it will converge and to what.

## 8.5 Prediction step: using the learned $\theta$

Once we have learned $\hat{\theta}$ using training data, given a new feature vector, $\boldsymbol{x}$, we can predict $y$ as

$$\hat{y} = [1, \boldsymbol{x}^\top]^\top \hat{\theta}$$

## 8.6 Understanding squared loss: Maximum Likelihood Estimation under i.i.d. Gaussian model **

The above can be motivated as Maximum Likelihood Estimation under the following probabilistic model for the data $\{y, \boldsymbol{x}\}$

$$p(y; \boldsymbol{x}, \theta) = \frac{1}{\sqrt{2\pi}\sigma_e} \exp((y - \boldsymbol{x}^\top \theta)^2 / 2\sigma_e^2)$$

Thus given $m$ training data points all of which are independent identically distributed (i.i.d.), we have the following model: $y^{(i)} = h_\theta(\boldsymbol{x}^{(i)}) + e^{(i)}$, $i = 1, 2, \ldots, m$ with $e^{(i)}$ standard Gaussian. Here the randomness is only in the noise $e$. In other words

$$\boldsymbol{y} \sim \mathcal{N}(\boldsymbol{X}\theta, \sigma_e^2 \boldsymbol{I})$$

or that

$$p(\boldsymbol{y}; \boldsymbol{X}, \theta) = \frac{1}{(\sqrt{2\pi}\sigma_e)^m} \exp\left(-\frac{\|\boldsymbol{y} - \boldsymbol{X}\theta\|^2}{2\sigma_e^2}\right)$$

*Maximum Likelihood Estimation (MLE)* means: find the value of $\theta$ that maximizes $p(\boldsymbol{y}; \boldsymbol{X}, \theta)$ (above probability density function). This is equivalent to minimizing the negative logarithm of the above PDF. It can be seen that this is equivalent to minimizing $J(\theta)$ given above and repeated here

$$J(\theta) = \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2$$

## 8.7 Least Squares, Regularized LS, Recursive LS and more

See Sec. 6.

# 9 Supervised Learning: Logistic Regression

In this case, $\boldsymbol{x}$ is still a real-valued $d \times 1$ vector but now $y$ is a scalar. The goal is classification. For two-class case, $y$ takes values either 0 or 1.

## 9.1 Two-class Model and its probabilistic interpretation

This assumes that $\Pr(y = 1; x, \theta) = h_\theta(\boldsymbol{x})$ with

$$h_\theta(\boldsymbol{x}) = g(\theta^T \boldsymbol{x}), \ \ g(z) := \frac{1}{1 + e^{-z}}$$

$g(.)$ is called the sigmoid function, it takes values between zero and one for all values of z. Thus, it can be used to model a probability. Said another way,

$$p(y; \boldsymbol{x}, \theta) = h_\theta(\boldsymbol{x})^y (1 - h_\theta(\boldsymbol{x}))^{1-y}$$

The prediction is

$$\hat{y} = \arg \max_{y=0,1} p(y; \boldsymbol{x}, \theta)$$

Thus

$$\hat{y} = 1 \text{ if } h_\theta(\boldsymbol{x}) > 1 - h_\theta(\boldsymbol{x}),$$

and $\hat{y} = 0$ otherwise. This simplifies to

$$\hat{y} = 1 \text{ if } \theta^T \boldsymbol{x} > 0$$

and $\hat{y} = 0$ otherwise.

### 9.1.1 Use of bias term in Logistic regression

Introduce the bias term exactly as we did in case of linear regression Make both $\theta$ and $\boldsymbol{x}$ $n = (d+1)$ length vectors; set the first entry of $\boldsymbol{x}$ equal to 1. The first entry of $\theta$ is then the bias term.

### 9.1.2 Learning $\theta$: Maximum Likelihood Estimation

Again define $\boldsymbol{y}$ and $\boldsymbol{X}$ as before from training data

Use Maximum Likelihood Estimation again: assume i.i.d. training data points $y^{(i)}$ (recall that this was assumed also in linear regression – it was imposed by letting the $e^{(i)}$'s be i.i.d.).

Thus, we minimize the negative log likelihood,

$$J(\theta) := -\log p(\boldsymbol{y}|\boldsymbol{X}; \theta) = -\log \left( \prod_{i=1}^m p(\boldsymbol{y}_i; \boldsymbol{x}^{(i)}, \theta) \right) = -\log \left( \prod_{i=1}^m h_\theta(\boldsymbol{x}^{(i)})^{y^{(i)}} (1 - h_\theta(\boldsymbol{x}^{(i)}))^{1-y^{(i)}} \right)$$

We can simplify this a lot as follows.

$$J(\theta) = -\sum_{i=1}^m \log \left( \left( \frac{1}{1 + \exp(\theta^\top \boldsymbol{x}^{(i)})} \right)^{y^{(i)}} \left( \frac{\exp(\theta^\top \boldsymbol{x}^{(i)})}{1 + \exp(\theta^\top \boldsymbol{x}^{(i)})} \right)^{1-y^{(i)}} \right)$$

$$= -\sum_{i=1}^m \log \left( \left( \frac{1}{1 + \exp(\theta^\top \boldsymbol{x}^{(i)})} \right)^{y^{(i)}} \left( \frac{1}{1 + \exp(-\theta^\top \boldsymbol{x}^{(i)})} \right)^{1-y^{(i)}} \right)$$

The only difference between the first and second terms inside $\log(.)$ is the sign of $\theta^\top \boldsymbol{x}^{(i)}$ in the denominator and the power it is raised to. One is raised to the power $y^{(i)}$, the other is raised to the power $(1 - y^{(i)})$. Here

$y^{(i)}$ is either 0 or 1. So, $2y^{(i)} - 1$ is either $-1$ (when $y^{(i)} = 0$) or $+1$ (when $y^{(i)} = 1$). Thus, when $y^{(i)}$ takes only values 0 or 1,

$$\left(\frac{1}{1 + \exp(\theta^\top \boldsymbol{x}^{(i)})}\right)^{y^{(i)}} \left(\frac{1}{1 + \exp(-\theta^\top \boldsymbol{x}^{(i)})}\right)^{1 - y^{(i)}} = \frac{1}{1 + \exp\left((\theta^\top \boldsymbol{x}^{(i)})(2y^{(i)} - 1)\right)}$$

The reason the last equality is true is because when $y^{(i)} = 1$, $2y^{(i)} - 1 = 1$, but when $y^{(i)} = 0$, $2y^{(i)} - 1 = -1$.

We now simply get

$$J(\theta) = -\sum_{i=1}^{m} \log\left(\frac{1}{1 + \exp\left((\theta^\top \boldsymbol{x}^{(i)})(2y^{(i)} - 1)\right)}\right) = \sum_{i=1}^{m} \log\left(1 + \exp\left((\theta^\top \boldsymbol{x}^{(i)})(2y^{(i)} - 1)\right)\right)$$

We can find $\theta$ by minimizing $J(\theta)$ by GD. The gradient is specified in the summary section given below. To compute the gradient, use the fact that the derivative of $\log g(z) = \log \frac{1}{1 + \exp(-z)} = -\log(1 + \exp(-z))$ w.r.t. $z$ is $\frac{\exp(-z)}{1 + \exp(-z)} = g(z)$.

It is possible to show that $J(\theta)$ is convex. Argument: weighted sum of convex functions is convex when the weights are positive, here the weights are just 1; $(\theta^\top \boldsymbol{x})(2y - 1)$ is an affine function of $\theta$ (and hence is both convex and concave; the logistic function $\log(1 + \exp(-z))$ can be shown be convex (see next line); composition of a convex function and an affine function is convex. To show $\log(1 + \exp(-z))$ is convex, since it is twice differentiable, we can compute the second derivative and show that it is $-\frac{1}{1 + \exp(-z))^2} < 0$ for any $z$, thus it is convex everywhere.

## 9.2 Two class Logistic Regression summary

Let $\theta$ be a $n = d + 1$ length vector with the first entry corresponding to the bias term. Replace $\boldsymbol{x}$ by $[1, \boldsymbol{x}^\top]^\top$ (this it makes it a $n = d + 1$ length vector

$$h_\theta(\boldsymbol{x}) = g(\theta^T \boldsymbol{x}), \ g(z) := \frac{1}{1 + e^{-z}}$$

$g(z) = \frac{1}{1 + \exp(-z)} = \frac{\exp(z)}{\exp(z) + 1}$ is called the sigmoid function, $0 \le g(z) \le 1$ for all $z$.

- Let $\theta$ be a $n = d + 1$ length vector with the first entry corresponding to the bias term. Replace $\boldsymbol{x}$ by $[1, \boldsymbol{x}^\top]^\top$ (this it makes it a $n = d + 1$ length vector

- Given training data $\boldsymbol{X}, \boldsymbol{y}$, use cost func as negative log likelihood. Assume independence of different data pairs.

$$J(\theta) := -\log p(\boldsymbol{y}|\boldsymbol{X}; \theta) = -\log\left(\prod_{i=1}^{m} p(\boldsymbol{y}^{(i)}; \boldsymbol{x}^{(i)}, \theta)\right) = -\log\left(\prod_{i=1}^{m} h_\theta(\boldsymbol{x}^{(i)})^{y^{(i)}} (1 - h_\theta(\boldsymbol{x}^{(i)}))^{1 - y^{(i)}}\right)$$

- We can estimate $\theta$ by minimizing $J(\theta)$ by GD. After simplification the gradient has the following expression:

$$\nabla J(\theta) = \sum_{i=1}^{m} (g(\theta^\top \boldsymbol{x}^{(i)}) - \boldsymbol{y}^{(i)}) \boldsymbol{x}^{(i)}$$

  - Use algorithm given in Sec. 5.1 with above gradient

- Classification: Given a query, $\boldsymbol{x}$, predict $y$ as

$$\hat{y} = 1 \text{ if } \theta^\top \boldsymbol{x} > 0, \ \hat{y} = 0 \text{ otherwise}$$

- Note: do not use the above model for simulating data on which to test the algorithm: for simulating data, better to use the GDA generative model.

### 9.2.1 Multi-class Log Reg

See Appendix

# 10 Supervised Learning: Generative Learning (a.k.a. Bayesian models/learning)

For logistic or linear regression we just assumed a probabilistic model on how $y$ is generated from $\boldsymbol{x}$, with $\boldsymbol{x}$ being deterministic.

In Generative Learning, we assume a "generative model": we first put a prior probabilistic model on $y$, and then assume a probabilistic model on how $\boldsymbol{x}$ was generated from $y$. We then compute the probability (or probability density function in case $y$ is real-valued) of $y$ taking a certain value given $\boldsymbol{x}$ using Bayes rule. Mathematically, we assume that we are given

$$p(\boldsymbol{x}|y;\theta), p(y)$$

and we use these to obtain the prediction as follows

$$\hat{y} = \arg\max_y p(y|\boldsymbol{x};\theta) := \arg\max_y \frac{p(\boldsymbol{x}|y;\theta)p(y;\theta)}{p(\boldsymbol{x};\theta)} = \arg\max_y p(\boldsymbol{x}|y;\theta)p(y;\theta)$$

This use of Bayes rule is called **Maximum A Posteriori (MAP)** detection or estimation in other literature. The overall approach is often called Bayesian modeling or physics-based modeling.

**Naive Bayes assumption:** This is often used to simplify the problem (reduce number of parameters). It says that conditioned on the class label $y$, the different features in $\boldsymbol{x}$ are mutually independent, i.e.,

$$p(\boldsymbol{x}|y;\theta) = \prod_{j=1}^{n} p(x_j|y;\theta).$$

With this, the expression for the posterior probability of the class label $y$ given $\boldsymbol{x}$ becomes

$$p(y|\boldsymbol{x};\theta) = \frac{p(\boldsymbol{x}|y;\theta)p(y;\theta)}{p(\boldsymbol{x};\theta)} = \frac{\prod_{j=1}^{n} p(x_j|y;\theta)p(y;\theta)}{\sum_{y'} \prod_{j=1}^{n} p(x_j|y';\theta)p(y')}$$

## 10.1 Learning $\theta$: Maximum Likelihood Estimation

Estimate $\theta$: define $\boldsymbol{y}$, $\boldsymbol{X}$ as before from training data. Also assume training data points are independent: $\{\boldsymbol{x}^{(i)}, y^{(i)}\}$ are mutually independent for different $i$. Define the cost function

$$J(\theta) := \Pr(\boldsymbol{y}, \boldsymbol{X};\theta) = \prod_{i=1}^{m} p(y^{(i)}, \boldsymbol{x}^{(i)};\theta)$$

or, usually its logarithm, and maximize it over $\theta$.

## 10.2 Generative Learning: Gaussian Discriminant Analysis (GDA)

This is one type of generative model and learning algorithm for the setting $\boldsymbol{x}$ real-valued $n \times 1$ vector and $y$ binary scalar. Thus $y$ can take two values 0 or 1. It assumes $\boldsymbol{x}$ is Gaussian given $y$ and $y$ itself is Bernoulli, i.e.,

$$p(\boldsymbol{x}|y;\theta) = \mathcal{N}(\boldsymbol{x};\mu_y, \boldsymbol{\Sigma}_y), \ p(y) = \phi^y(1-\phi)^{1-y}$$

Notice in this case $\theta = \{\mu_0, \mu_1, \boldsymbol{\Sigma}_0, \boldsymbol{\Sigma}_1, \phi\}$. Parameters are still learnt by MLE

$$\max_\theta J(\theta) := \Pr(\boldsymbol{y}, \boldsymbol{X};\theta) = \prod_{i=1}^{m} p(y^{(i)}, \boldsymbol{x}^{(i)};\theta) \text{ s.t. } 0 \le \phi \le 1$$

Notice that, without extra assumptions, we have $2n + 2n^2 + 1$ parameters. Training data are each $n$ length vectors $\boldsymbol{x}^{(i)}$, thus we can say we have $mn$ training data scalars. We need $mn$ significantly larger than $2n + 2n^2 + 1$ for training/learning to be accurate. We will need $m$ growing at least linearly with $n$ to be able to learn anything useful.

But the point of Bayesian (generative) modeling is that we should be able to use a smaller $m$ and still train well.

When enough training data is not available, we need to simplify our model so that there are fewer parameters. As explained later, this will increase model bias, but will reduce the variance in parameter estimation.

A common model simplification is to assume that the different entries of each $\boldsymbol{x}^{(i)}$ are independent conditioned on the class label $y^{(i)}$. This is called the *Naive Bayes assumption*.

### 10.2.1 Gaussian Discriminant Analysis (GDA) with Naive Bayes assumption and equal covariances

A common model simplification is to assume that the different entries of each $\boldsymbol{x}^{(i)}$ are independent conditioned on the class label $y^{(i)}$. This is called the *Naive Bayes assumption*.

In the Gaussian case, this translates to assuming that $\Sigma_0, \Sigma_1$ are *diagonal*. With the diagonal assumption, we now have only $2n + 2n + 1$ parameters which is much more manageable. A second commonly used simplification is to assume the same covariance under both classes, i.e., that $\Sigma_0 = \Sigma_1 = \Sigma$ and $\Sigma$ is diagonal. With this assumption too, we have the following simpler model

$$\prod_{i=1}^{m} \left( \left( \prod_{j=1}^{n} \mathcal{N}(\boldsymbol{x}_j; (\mu_{y_i})_j, \sigma_j^2) \right) \cdot \phi^{y_i} (1 - \phi)^{1 - y_i} \right)$$

Under the above assumption, the Max Likelihood Estimates (MLE) of the model parameters, $\theta := \{\phi, \mu_0, \mu_1, \Sigma\}$ i.e., the value of the model parameters that solve

$$\arg \max_{\theta} J(\theta), \; J(\theta) := \Pr(\boldsymbol{y}, \boldsymbol{X}; \theta) = \prod_{i=1}^{m} p(y^{(i)}, \boldsymbol{x}^{(i)}; \theta) \text{ s.t. } 0 \le \phi \le 1$$

are computed as follows:

$$\hat{\phi} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{1}(y^{(i)} = 1)$$

$$\hat{\mu}_0 = \frac{\sum_{i=1}^{m} \mathbf{1}(y^{(i)} = 0) \boldsymbol{x}^{(i)}}{\sum_{i=1}^{m} \mathbf{1}(y^{(i)} = 0)}$$

$$\hat{\mu}_1 = \frac{\sum_{i=1}^{m} \mathbf{1}(y^{(i)} = 1) \boldsymbol{x}^{(i)}}{\sum_{i=1}^{m} \mathbf{1}(y^{(i)} = 1)}$$

$$\hat{\sigma_j^2} = \frac{1}{m} \sum_{i=1}^{m} (\boldsymbol{x}^{(i)} - \mu_{y^{(i)}})_j^2, \; j = 1, 2, \ldots, n$$

while setting all non-diagonal entries of $\hat{\Sigma}$ to be zero. Here $\mathbf{1}$ denotes the *indicator function* of the statement in paranthesis. Thus, $\mathbf{1}(y^{(i)} = 0)$ equals one if $y^{(i)} = 0$ and it equals zero otherwise.

Notice that the above is equivalent to learning the parameters for *each feature* independently, it can also

be rewritten as follows: for each $j = 1, 2, \ldots, n$, compute

$$(\hat{\mu}_0)_j = \frac{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 0)(\boldsymbol{x}^{(i)})_j}{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 0)}$$

$$(\hat{\mu}_1)_j = \frac{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 1)(\boldsymbol{x}^{(i)})_j}{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 1)}$$

$$\hat{\sigma}_j^2 = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{x}^{(i)} - \mu_{y^{(i)}})_j^2, \ \ j = 1, 2, \ldots, n$$

In this Naive Bayes' setting, the selecting of $\hat{y}$ to solve

$$\hat{y} = \arg\max_y p(\boldsymbol{x}|y; \hat{\theta}) p(y; \hat{\theta})$$

simplifies to the following test (just simplify the expressions)

$$\hat{y} = 1 \text{ if } \sum_{j=1}^n \frac{(\boldsymbol{x} - \hat{\mu}_1)_j^2}{\hat{\sigma}_j^2} - \sum_j \frac{(\boldsymbol{x} - \hat{\mu}_0)_j^2}{\hat{\sigma}_j^2} < thresh$$

with *thresh* depending on $\hat{\phi}$. If $\hat{\phi} \approx 0.5$ then *thresh* $= 0$.

### 10.2.2  Notes

1. In the Gaussian case, we get a simple closed form as above. In many cases of Bayesian modeling also, this is possible. In general when this is not possible, *do not* ever work directly with probabilities. *Always work with logarithms of the probabilities.* Otherwise you will run into numerical problems while coding. Remember: under the naive Bayes assumption, $p(\boldsymbol{x}|y; \hat{\theta}) = \prod_{j=1}^n p(\boldsymbol{x}_j|y; \hat{\theta})$ if $p(.)$ is a PMF, it is a product of $n$ real numbers all less than one. Even if $p(.)$ is a PDF (which could be more than one), very often it is actually less than one only. Product n numbers less than one can become very small very soon.

2. When working with real data, e.g., images, there may be a certain region that is black in *all* the images. For these pixels the estimate of the variance will be zero meaning $\sigma_j^{-1} = \infty$ resulting in code bugs.

   (a) Fix 1: do not use these directions.

   (b) Fix 2: Better fix to deal with similar issues where in some directions variance is very small (may not be zero): use PCA to reduce the dimensionality of the data. For classification there is no need to use all $n$ features.

   (c) Fix 3: add a small value to replace the zeros: the added value should be *much smaller* than any of the *important directions' variances*. This is easy in the zero/nonzero case but in practice ill-conditioning of $\Sigma$ may make this hard to do. Use of PCA and a smaller dimension is thus a much better fix.

## 10.3  Generative Learning: Spam Filter – an example of Discrete-valued or Categorical Features

In applications such as spam email detector (or filter) design, one typically models $\boldsymbol{x}$ as a discrete-valued vector given $y$.

$y = 0$ means the email is not-spam, $y = 1$ means it is spam.

### 10.3.1  Simple Spam Filter: entries of x are binary

In the simplest version, $\boldsymbol{x}$ is an $n \times 1$ binary vector with $n$ being equal to the size of the English dictionary. We say $\boldsymbol{x}_j = 1$ if the $j$-th dictionary word is in the email and $\boldsymbol{x}_j = 0$ otherwise. This means $n$ is really large. Also it is not counting how many times a word occurred.

Consider first the simplest model where $\boldsymbol{x}$ is a binary vector. As before, we specify

$$p(\boldsymbol{x}|y; \theta), \ p(y)$$

and we predict

$$\hat{y} = \arg\max_{y} p(\boldsymbol{x}|y; \theta)p(y; \theta).$$

Notice now that $\boldsymbol{x}$ can take a total of $2^n$ possible values. We also need to specify the prior $p(y)$. Thus, in this most general case, the number of parameters equals $2^n + 1$.

Training data are each $n$ length vectors $\boldsymbol{x}^{(i)}$, thus we can say we have $mn$ training data scalars. We need $mn$ significantly larger than $2^n + 1$ for training/learning to be accurate. Here we will need $m$ to grow linearly with $2^{n-1}$: this can be very large and is not practical.

### 10.3.2  Naive Bayes assumption

In both the above examples and especially the second one, the required $m$ can be very large for accurate training/learning. Thus we add a further modeling assumption called "Naive Bayes" in ML literature. Others would call it *"conditional independence" of different entries of a feature vector (the different $\boldsymbol{x}_j$'s, $j = 1, 2, \ldots, n$) given $y$.* Mathematically, we are assuming

$$p(\boldsymbol{x}|y; \theta) = \prod_{j=1}^{n} p(\boldsymbol{x}_j|y; \theta)$$

This may not be a very realistic assumption, but it significantly reduces the number of parameters required by the model.

In the Gaussian case, this implies that $\boldsymbol{\Sigma}_0, \boldsymbol{\Sigma}_1$ are **diagonal matrices**. Thus, the number of parameters becomes $2n + 2n + 1$ which is much more tractable. In the spam filter case, this means we have $n + 1$ parameters.

So now the number of training samples $m$ does not even need to grow with $n$.

### 10.3.3  Simple Spam filter with Naive Bayes

$\boldsymbol{x}$ is a binary vector, $y$ is a scalar. With using Naive Bayes, in the Simple Spam Filter case, we can now define

$$\psi_{j,y} := p(\boldsymbol{x}_j = 1|y), j = 1, 2, \ldots, n; \ and \ \phi := p(y = 1)$$

With this we have just $n + 1$ parameters to learn instead of $2^n + 1$.

We can again learn the parameters by MLE:

$$\max_{\theta} J(\theta) := \prod_{i=1}^{m} p(y^{(i)}, \boldsymbol{x}^{(i)}; \theta) = \prod_{i=1}^{m} \prod_{j=1}^{n} p(\boldsymbol{x}_j^{(i)}|y^{(i)}; \theta)p(y^{(i)}; \theta) = \prod_{i=1}^{m} \phi^{y^{(i)}}(1-\phi)^{(1-y^{(i)})} \left( \prod_{j=1}^{n} \psi_{j,y}^{\boldsymbol{x}_j^{(i)}} (1-\psi_{j,y})^{(1-\boldsymbol{x}_j^{(i)})} \right)$$

s.t. constraints that

$$0 \le \psi_{j,y} \le 1, \ 0 \le \phi \le 1$$

Can again get closed form simple expressions for the MLE:

$$\hat{\psi}_{j,0} =$$

count the number of training data points for which $y^{(i)} = 0$ and $\boldsymbol{x}_j^{(i)} = 1$ and divide by the total number of training data points for which $y^{(i)} = 0$.

Similarly

$$\hat{\psi}_{j,1} =$$

count the number of training data points for which $y^{(i)} = 1$ and $\boldsymbol{x}^{(i)} = 1$ and divide by the total number of training data points for which $y^{(i)} = 1$. and

$$\hat{\phi} =$$

count the number of training data points for which $y^{(i)} = 1$ and divide by $m$

Use of $n$ to be dictionary size makes it extremely large causing the algorithm to be very slow. Also, instead of letting $n$ be the size of the English dictionary, we can let $n$ be the size of the vocabulary (set of all words in all training data). But this has the disadvantage that it does not tell you how to deal with an unseen word. Options include (i) ignore unseen words (can be problematic in case the unseen words are the only reason an email is obviously spam); or (ii) in the model, assume a small nonzero probability for an unseen word (this means: increase the vocabulary size by 1, this probability cannot be learned easily, you just have to make up a "reasonable" value for it).

# 11 Supervised Learning: Linear Classifiers and Support Vector Machines (SVMs)

## 11.1 Linear Classifiers **

Both logistic regression and GDA with Naive Bayes and equal covariances result in a linear classifier, i.e., one can simplify the classification rule in both cases to get the following:

$$\hat{y} = 1 \ \text{ if } \boldsymbol{w}^T \boldsymbol{x} + b > 0$$

and equals 0 otherwise.

Proof for GDA:

GDA decision rule is: $\hat{y} = 1$ if

$$(\boldsymbol{x} - \mu_1)^T \Sigma_1^{-1} (\boldsymbol{x} - \mu_1) < (\boldsymbol{x} - \mu_0)^T \Sigma_0^{-1} (\boldsymbol{x} - \mu_0)$$

With $\Sigma_1 = \Sigma_0 = \Sigma$ and naive Bayes ($\Sigma$ is diagonal), this simplifies to checking if

$$\sum_j \frac{(\boldsymbol{x} - \mu_1)_j^2}{\sigma_j^2} < \sum_j \frac{(\boldsymbol{x} - \mu_0)_j^2}{\sigma_j^2}$$

Simplifying the above, we equivalently need

$$\boldsymbol{w}^T \boldsymbol{x} + b > 0, \quad \text{with} \quad w_j = \frac{(\mu_1 - \mu_0)_j}{\sigma_j^2}, \quad b = 0.5 \sum_j \frac{(\mu_1)_j^2 - (\mu_0)_j^2}{\sigma_j^2}$$

In fact we can also get an expression for $\boldsymbol{w}$ and $b$ even if we just have $\Sigma_1 = \Sigma_0 = \Sigma$

$$\boldsymbol{w} = \Sigma^{-1}(\mu_1 - \mu_0), \quad b = 0.5(\mu_1^T \Sigma^{-1} \mu_1 - \mu_0^T \Sigma^{-1} \mu_0).$$

## 11.2 Support Vector Machines (SVMs): Motivation and main idea

Both the classifiers we talked about so far are linear classifiers as explained above. Consider logistic regression. The classification decision would be much more reliable if $\theta^T \boldsymbol{x}$ were either much larger or much smaller than zero.

The idea of SVM is this: we do not assume any data model here. Instead we try to look for the "separating hyperplane", equivalently, a vector $\theta$ so that the "margin" from the decision boundary is maximized for *all* training data points. Visual explanation in class or see cs229-notes-3.

### 11.2.1 Notation change

Instead of a single vector $\theta$ with the first entry used for the bias term, in case of SVMs, we use a weight vector $\boldsymbol{w}$ which is the same length as the data and a scalar $b$. Also, instead of labeling the two classes as 0 and 1, we label them is $-1$ and $+1$ because this simplifies some of the writing.

Margin: the distance of a data point from the separating hyperplane. Margin for the $i$-th training data point is computed as

$$y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b)$$

### 11.2.2 Using SVM for classification

Suppose first that the optimal choice of $\boldsymbol{w}, b$ is available. Then we classify as follows

$$\hat{y} = sign(\boldsymbol{w}^T \boldsymbol{x} + b)$$

If the term $> 0$, then the class is $+1$ else the class is $-1$.

### 11.2.3 Goal

The goal in case of SVMs is to find $\boldsymbol{w}, \boldsymbol{b}$ that maximize the worst-case margin defined by

$$\min_{i=1,2,\ldots,m} y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b)$$

By multiplying $\boldsymbol{w}, b$ by a scalar we could keep increasing the margin, but that will not improve classification. Thus, we need to impose the constraint that $\|\boldsymbol{w}\|_2 = 1$ or equivalently, replace $\boldsymbol{w}$ by $\boldsymbol{w}/\|\boldsymbol{w}\|$.

## 11.3 Simplifications to obtain a convex optimization problem **

Writing a slightly different way, we need to solve

$$\max_{\boldsymbol{w},b} \gamma \text{ s.t. } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq \gamma, \ i = 1, 2, \ldots, m \text{ and } \|\boldsymbol{w}\|_2 = 1$$

In the above, notice that we could scale everything by any scalar, say multiply through by $1/(\|\boldsymbol{w}\|)$, and the problem will not change, i.e.,

$$\max_{\boldsymbol{w},b} \frac{\gamma}{\|\boldsymbol{w}\|} \text{ s.t. } y^{(i)}(\frac{1}{\|\boldsymbol{w}\|}\boldsymbol{w}^T \boldsymbol{x}^{(i)} + \frac{1}{\|\boldsymbol{w}\|}b) \geq \frac{\gamma}{\|\boldsymbol{w}\|}, \ i = 1, 2, \ldots, m \text{ and } \|\boldsymbol{w}\|_2 = 1$$

Let $\tilde{b} = b/\|\boldsymbol{w}\|$, and $\tilde{\boldsymbol{w}} = \boldsymbol{w}/\|\boldsymbol{w}\|$. Then we can optimize over

$$\max_{\boldsymbol{w},\tilde{b}} \frac{\gamma}{\|\boldsymbol{w}\|} \text{ s.t. } y^{(i)}(\tilde{\boldsymbol{w}}^T \boldsymbol{x}^{(i)} + \tilde{b}) \geq \frac{\gamma}{\|\boldsymbol{w}\|}, \ i = 1, 2, \ldots, m \text{ and } \|\boldsymbol{w}\|_2 = 1$$

The above is equivalent to dividing everything by $\|\boldsymbol{w}\|_2$. Doing this gives

$$\max_{\boldsymbol{w},b} \frac{\gamma}{\|\boldsymbol{w}\|_2} \text{ subject to } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq \gamma, \ i = 1, 2, \ldots, m$$

Writing the above, we need to solve

$$\max_{\boldsymbol{w},b} \gamma \text{ s.t. } \frac{y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b)}{\|\boldsymbol{w}\|} \geq \gamma, \ i = 1, 2, \ldots, m$$

This is not a convex optimization problem yet. So we try to simplify further. This is the same as

$$\max_{\boldsymbol{w},b} \gamma \text{ s.t. } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq \gamma\|\boldsymbol{w}\|, \ i = 1, 2, \ldots, m$$

Let $\tilde{\gamma} = \gamma \|\boldsymbol{w}\|$, then we have

$$\max_{\boldsymbol{w},b} \frac{\tilde{\gamma}}{\|\boldsymbol{w}\|} \text{ subject to } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq \tilde{\gamma}, \ i = 1, 2, \ldots, m$$

This is still not convex. Another point to notice is this: we could fix the value of the margin $\gamma$ to 1, and nothing will change. This gives

$$\max_{\boldsymbol{w},b} \frac{1}{\|\boldsymbol{w}\|_2} \text{ subject to } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq 1, \ i = 1, 2, \ldots, m$$

This is not convex either but now a simple reformulation gives us a convex problem. Maximizing $1/\|\boldsymbol{w}\|_2$ is equivalent to minimizing $\|\boldsymbol{w}\|_2$ which is the same as minimizing $\|\boldsymbol{w}\|_2^2$. This gives

$$\min_{\boldsymbol{w},b} \|\boldsymbol{w}\|_2^2 \text{ subject to } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq 1, \ i = 1, 2, \ldots, m$$

This is now a convex optimization problem since the cost is a convex function and the inequality constraints are linear. In particular it is what is called a Quadratic Program or QP.

## 11.4   Final primal problem

$$\min_{\boldsymbol{w},b} \|\boldsymbol{w}\|_2^2 \text{ subject to } y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq 1, \ i = 1, 2, \ldots, m$$

This is now a convex optimization problem since the cost is a convex function and the inequality constraints are linear. In particular it is what is called a Quadratic Program or QP.

## 11.5   Simplification using duality: faster optimization when $m \ll n$

By using *Lagrange duality*, it is possible to show that we can also compute the optimal $\boldsymbol{w}, b$ as follows.

1. Solve the following "dual problem": optimize over the Langrange multipliers $\alpha_i$

$$\max_{\alpha} \left( \sum_{i=1}^{m} \alpha_i - 0.5 \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)} \rangle \right) \text{ s.t. } \sum_{i=1}^{m} \alpha_i y^{(i)} = 0, \ \alpha_i \geq 0, \ i = 1, 2, \ldots, m$$

   or equivalently solve the following minimization problem: note change of sign

$$\min_{\alpha} \left( - \sum_{i=1}^{m} \alpha_i + 0.5 \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)} \rangle \right) \text{ s.t. } \sum_{i=1}^{m} \alpha_i y^{(i)} = 0, \ \alpha_i \geq 0, \ i = 1, 2, \ldots, m$$

   Get $\hat{\alpha}$ as output

   (here $\langle \boldsymbol{x}_1, \boldsymbol{x}_2 \rangle = \boldsymbol{x}_1^T \boldsymbol{x}_2$ is the inner product)

2. Obtain

$$\hat{\boldsymbol{w}} = \sum_{i=1}^{m} \hat{\alpha}_i y^{(i)} \boldsymbol{x}^{(i)} \quad \text{(this never needs to be computed, see notes below)}$$

$$\hat{b} = -0.5 \left( \max_{i:y^{(i)}=-1} \hat{\boldsymbol{w}}^T \boldsymbol{x}^{(i)} + \min_{i:y^{(i)}=1} \hat{\boldsymbol{w}}^T \boldsymbol{x}^{(i)} \right)$$

3. Classification step:

$$\hat{y} = sign(\hat{\boldsymbol{w}}^T \boldsymbol{x} + \hat{b}) = sign(\sum_i \alpha_i y^{(i)} \langle \boldsymbol{x}^{(i)}, \boldsymbol{x} \rangle + \hat{b})$$

where

$$\hat{b} = -0.5 \left( \max_{i:y^{(i)}=-1} \sum_{j=1}^m \hat{\alpha}_j y^{(j)} \langle \boldsymbol{x}^{(j)}, \boldsymbol{x}^{(i)} \rangle + \min_{i:y^{(i)}=1} \sum_{j=1}^m \hat{\alpha}_j y^{(j)} \langle \boldsymbol{x}^{(j)}, \boldsymbol{x}^{(i)} \rangle \right)$$

(note: this also only uses inner products: for a query feature $\boldsymbol{x}$ – use for later in Kernel SVM)

Notice: In the above dual program

- $\alpha$ is a vector of length $m$: for classification usually $m \ll n$ suffices. Thus the dual opt problem is computationally much cheaper to solve

- The dependence on feature vectors is only through inner products: this is true for training (for obtaining $\hat{\alpha}$). This is also true for testing: the weight vector $\hat{\boldsymbol{w}}$ actually never needs to be computed. We only need to compute the $\hat{\alpha}$ vector: we can then use inner products between the query $\boldsymbol{x}$ and the training data features and the $\alpha$ vector to compute the output $\hat{y}$ for a given $\boldsymbol{x}$

The above is useful in two settings: (i) if $n \gg m$, i.e. the original data lies in a much higher dimensional space compared to the available number of data points (this often happens in classification problems), then the above dual is much less expensive to solve. (ii) Notice that in the above problem, all dependence on the feature vectors $\boldsymbol{x}^{(i)}$ is through the inner product $\langle \boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)} \rangle$. This means we could consider higher dimensional feature vectors, i.e., convert $\boldsymbol{x}$ to $\phi(\boldsymbol{x})$ by considering features of the form $\boldsymbol{x}_j^2$, $\boldsymbol{x}_j \boldsymbol{x}_k$.

## 11.6 SVM Dual Problem and CVXOPT Mapping

We consider the dual form of the hard-margin SVM. Given training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$, where $\mathbf{x}_i \in \mathbb{R}^n$ and $y_i \in \{-1, +1\}$, the dual optimization problem is:

$$\begin{aligned}
\max_{\boldsymbol{\alpha} \in \mathbb{R}^m} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\
\text{subject to} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\
& \alpha_i \geq 0 \quad \text{for all } i.
\end{aligned} \tag{2}$$

This can be cast into the standard form of a quadratic programming (QP) problem:

$$\begin{aligned}
\min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \boldsymbol{\alpha}^\top P \boldsymbol{\alpha} + q^\top \boldsymbol{\alpha} \\
\text{subject to} \quad & G \boldsymbol{\alpha} \preceq h, \\
& A \boldsymbol{\alpha} = b,
\end{aligned} \tag{3}$$

where $\boldsymbol{\alpha} \in \mathbb{R}^m$ is the vector of dual variables.

### 11.6.1 Term-by-Term Mapping to CVXOPT

### 11.6.2 Python Code Using `cvxopt`

```
from cvxopt import matrix, solvers
import numpy as np

# Assume X is m x n data matrix, y is m x 1 label vector (+1 or -1)
```

| SVM Notation | Description | CVXOPT Variable | Expression |
|---|---|---|---|
| $\boldsymbol{\alpha} \in \mathbb{R}^m$ | Dual variables | Decision variable | $\boldsymbol{\alpha}$ |
| $P \in \mathbb{R}^{m \times m}$ | Quadratic term | P | $P_{ij} = y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$ |
| $q \in \mathbb{R}^m$ | Linear term | q | $q_i = -1$ |
| $G \in \mathbb{R}^{m \times m}$ | Inequality matrix | G | $G = -I$ |
| $h \in \mathbb{R}^m$ | Inequality vector | h | $h = 0$ |
| $A \in \mathbb{R}^{1 \times m}$ | Equality constraint | A | $A = [y_1, \ldots, y_m]$ |
| $b \in \mathbb{R}$ | Equality value | b | $b = 0$ |

Table 1: Mapping between SVM dual terms and CVXOPT inputs

```python
K = np.dot(X, X.T)  # Gram matrix
P = matrix(np.outer(y, y) * K)
q = matrix(-np.ones(X.shape[0]))
G = matrix(-np.eye(X.shape[0]))
h = matrix(np.zeros(X.shape[0]))
A = matrix(y.reshape(1, -1).astype('double'))
b = matrix(np.zeros(1))

solution = solvers.qp(P, q, G, h, A, b)
alphas = np.array(solution['x'])
```

## 11.7   Soft margin SVMs

See page 19 of ML-cs229-noted-3

The SVM presented so far assumes data is linearly separable. But often data is not. Then we use the following

$$\min_{\boldsymbol{w}, b, \zeta} \|\boldsymbol{w}\|_2^2 + C_1 \sum_{i=1}^m \zeta^{(i)} \quad \text{s.t.} \quad y^{(i)}(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b) \geq 1 - \zeta^{(i)}, \ \zeta^{(i)} \geq 0, \ i = 1, 2, \ldots, m,$$

The above is imposing the following: find $\boldsymbol{w}, b$ so that some training data are misclassified but not too many.

Dual problem for the above

$$\max_{\alpha, \zeta} \left( \sum_{i=1}^m \alpha_i - 0.5 \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} \langle x^{(i)}, \boldsymbol{x}^{(j)} \rangle \right) \quad \text{s.t.} \quad \sum_{i=1}^m \alpha_i y^{(i)} = 0, \ \alpha_i \geq 0, \ \alpha_i \leq C_1, \ i = 1, 2, \ldots, m$$

1. If $\zeta^{(i)} > 1$, then the point $i$ is misclassified.

2. If $0 \leq \zeta^{(i)} < 1$, then the point is not misclassified but it is within the "margin" region (it is not being used to select the separating hyperplane).

3. $C = \infty$ or very large will result in $\zeta = \boldsymbol{0}$ being the solution, i.e., this then becomes the same as hard-margin SVM. No point can be misclassified.

4. $C = 0$ or very small means $\zeta$ can be anything. Then the best solution is $\boldsymbol{w} = \boldsymbol{0}$ and $\zeta$ anything, the constraint becomes $LHS > -\infty$ which is equivalent to no constraint. Thus $C = 0$ will return a completely wrong SVM.

5. Increase $C$: fewer misclassifications allowed.

6. Increase $C$: lower bias, higher variance

7. How to tune $C$ - vary in the range $10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1, 10, 100, \ldots$. Test on "tuning data". Split available data into training set, tuning set, test set.

## 11.8 Kernel SVMs: use the Dual program

Notice in the dual program that the dependence on feature vectors is only through inner products: this is true for training (for obtaining $\hat{\alpha}$). This is also true for testing: the weight vector $\hat{\boldsymbol{w}}$ actually never needs to be computed. We only need to compute the $\hat{\alpha}$ vector: we can then use inner products between the query $\boldsymbol{x}$ and the training data features and the $\alpha$ vector to compute the output $\hat{y}$ for a given $\boldsymbol{x}$.

The above implies we can replace the regular Euclidean inner product by any other inner product and our solution form will not change (in terms of coding, it would imply you need to code in a new function. Why would we do that? One reason:

- If the original data is not linearly separable, we try to map it to a higher dimensional space and hope that it is linearly separable in that space.

- How to do this? One example is use the original feature vector and pairwise products of the features, e.g. if $n = 2$, then use $\phi(\boldsymbol{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]^\top$.

- With a few constants changed, the inner product between $\phi(\boldsymbol{x}_a)$, $\phi(\boldsymbol{x}_b)$ for the above $\phi(.)$ mapping can be computed as
$$K(\boldsymbol{x}_a, \boldsymbol{x}_b) := < \phi(\boldsymbol{x}_a), \phi(\boldsymbol{x}_b) > = (\boldsymbol{x}_a^\top \boldsymbol{x}_b + c)^2$$

  The above is called the "kernel product"

- Notice computing the inner product using $(\boldsymbol{x}_a^\top \boldsymbol{x}_b + c)^2$ has time cost of order $n$, while computing the inner product directly for $\phi(\boldsymbol{x}_a), \phi(\boldsymbol{x}_b)$ has cost order $n_{higher} = n^2$.

- There are many other options for higher dimensional mappings and the corresponding Kernel products

- To add: more examples; Theorem for what can be a valid Kernel product, final conclusion

For many such high dimensional feature mappings, $\phi(\boldsymbol{x})$ is very expensive to compute. But just computing the inner product, defined as the "kernel product"

$$K_\phi(\boldsymbol{x}_1, \boldsymbol{x}_2) := \langle \phi(\boldsymbol{x}_1), \phi(\boldsymbol{x}_2) \rangle$$

is much less expensive. An example is $\phi(x)$ obtained as all pairwise products of entries of $\boldsymbol{x}$. For this, computing $\phi(\boldsymbol{x}^{(i)})$ takes order $n^2$ time, thus computing the new feature vector for all $i$ takes order $n^2 m$ time. But computing one kernel product can be done as

$$K_\phi(\boldsymbol{x}_1, \boldsymbol{x}_2) = (\boldsymbol{x}_1^T \boldsymbol{x}_2)^2$$

This only takes order $n$ time. There are $m(m+1)/2$ total products to compute. Thus, the time needed is of order $nm^2$. Since typically, $m \ll n$, this is much quicker.

It is also possible to define kernel products for cases where the actual feature mapping is infinite dimensional. Gaussian kernel is an example. This is defined by

$$K_\phi(\boldsymbol{x}_1, \boldsymbol{x}_2) = \exp\left( \frac{\|\boldsymbol{x}_1 - \boldsymbol{x}_2\|^2}{2\sigma^2} \right)$$

Kernel product: is basically some measure of similarity between two data points. Any useful measure of similarity can be used to define a "kernel" (kernel product), one may not even need to specify the underlying feature mapping.

A very large number of kernels can be defined. The purpose is for datasets which are not linearly separable in the original feature space, it is possible they are in a higher dimensional space.

### 11.8.1 Training and Classification using Kernel SVM

1. Training: Solve

$$\min_{\alpha} \left( -\sum_{i=1}^{m} \alpha_i + 0.5 \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y^{(i)} y^{(j)} K_{\phi}(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(j)}) \right) \ \text{s.t.} \ \sum_{i=1}^{m} \alpha_i y^{(i)} = 0, \ \alpha_i \geq 0, \ i = 1, 2, \ldots, m$$

Get $\hat{\alpha}$ as output.

(here $\langle \boldsymbol{x}_1, \boldsymbol{x}_2 \rangle = \boldsymbol{x}_1^T \boldsymbol{x}_2$ is the inner product)

2. Classify:

$$\hat{y} = sign(\sum_{i} \alpha_i y^{(i)} K_{\phi}(\boldsymbol{x}^{(i)}, \boldsymbol{x}) + \hat{b})$$

where

$$\hat{b} = -0.5 \left( \max_{i:y^{(i)}=-1} \sum_{j=1}^{m} \hat{\alpha}_j y^{(j)} K_{\phi}(\boldsymbol{x}^{(j)}, \boldsymbol{x}^{(i)}) + \min_{i:y^{(i)}=1} \sum_{j=1}^{m} \hat{\alpha}_j y^{(j)} K_{\phi}(\boldsymbol{x}^{(j)}, \boldsymbol{x}^{(i)}) \right)$$

### 11.8.2 Other ML problems

This "kernel trick" can be used for many other learning algorithms as well. Anytime all computation depends on inner products between the features, this can be used.

## 11.9 Deviation: Introduction to Langrange duality to understand how to derive the dual program **

To be updated later (to make this course 425/525).

See ML-cs229-notes3

To do also: add a section on basic optimization ideas - just enough to teach MLE derivations; and a section on MLE where we derive all the parameter estimates.

# 12 Unsupervised Learning: PCA

In unsupervised learning, there is no *labeled* training data to learn parameters from. So no "output" $y^{(i)}$ is available for "input" $\boldsymbol{x}^{(i)}$. PCA is an unsupervised learning technique that is used for dimension reduction. Given data vectors in $\Re^n$, if they approximately lie in a lower dimensional subspace, how do we find that subspace?

## 12.1 What is PCA

As before we assume that we are given $m$ data vectors (usually called feature vectors) $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \cdots, \boldsymbol{x}^{(m)}$ each in $\Re^n$, and these are stacked as the rows of a matrix

$$\boldsymbol{X} = \begin{bmatrix} -(\boldsymbol{x}^{(1)})^T- \\ -(\boldsymbol{x}^{(2)})^T- \\ \vdots \\ -(\boldsymbol{x}^{(m)})^T- \end{bmatrix} \tag{4}$$

Given a reduced dimension $r$, the goal is find an $r$ dimensional representation of the data vectors so that maximum information (variance) is retained.

## 12.2 Why PCA

Notice each data vector is $n$-length. $n$ may be very large can, e.g., can be equal to the image size if one uses all the pixels as "features". So we try to reduce the dimension to $r < n$. This is helpful in the following ways:

- Memory and Speed: Less storage needed to save the data. Analysis of the reduced-dimension dataset, e.g., regression or classification, will be faster;

- Noise reduction: if features are noisy, i.e., if $\boldsymbol{X}$ is such that $(1/m)\boldsymbol{X}\boldsymbol{X}^\top \approx \Sigma + \sigma_e^2 \boldsymbol{I}$.

- Removing feature vector entries (or directions) with zero or near zero variability in the training data

  - In some other cases, there may be zero variance along some directions and when the data co-variance matrix is computed, it ends up being rank deficient. Hence it cannot be inverted, but inversion is needed for example, for GDA based classification. This is actually the perfect application where PCA resolve the problem and may in fact help improve classification accuracy.

  - This is also useful for noise reduction when test data does not follow the same distribution as training data - and test data has noise even along these directions.

- Finally: PCA helps to get subspace coefficients which are uncorrelated. This is needed in some applications and is convenient in others. For example, for GDA things means that the covariance matrix of the new data vectors is diagonal.

  - For this, we need not reduce the dimension though, we can obtain uncorrelated variables in $n$ dimensions too by using the full SVD.

- PCA useful for linear regression only if (i) either the original set of features is very noisy with noise being i.i.d. across all directions; or (ii) if the features set is highly correlated (the resulting feature covariance is approx low rank). Even so, in most cases using $r = (m-1)$ will often be the best solution. The real question is: is $r = m/2$ almost as good?

## 12.3 Use of PCA and PCA vs Feature Normalization

PCA is a dimension reduction technique.

It is one way to do feature selection if the assumption that "features having the most variation are the most important ones" is true. In addition, of course PCA helps de-correlate the features.

A typical setting where this is useful is when the raw feature vector is all image pixels. There are too many pixels and pixels are highly correlated.

*PCA is not very useful if features are first variance normalized and then PCA is applied.*

### 12.3.1 What is Feature normalization

When ML features have vastly different ranges of values, e.g., in the house price example, one house feature can be number of bedrooms which takes values say between 1 to 5, while a second feature is the house areas in square feet, this could be a few 100 square feet. The two features have vastly different ranges.

If these were not normalized, then the matrix $\boldsymbol{X}$ would be ill-conditioned. This means that it will have a very high condition number. This, in turn, will imply inaccurate estimation if the closed form solution is used (due to numerical errors in inverting an ill conditioned $\boldsymbol{X}^\top \boldsymbol{X}$). If GD is used, the same can happen unless the step size is chosen to be very small. In this case, it will converge to the correct solution but will need too many iterations.

One solution is the normalize each feature: compute the $j$-th feature's mean $\mu_j$ and standard deviation $\sigma_j$; replace $\boldsymbol{x}_j$ by $(\boldsymbol{x}_j - \mu_j)/\sigma_j$.

### 12.3.2    Feature normalization vs PCA

Feature normalization is used in ML settings like the above where different features have different ranges because they are measured using very different measuring tools or sensors; and consequently, larger variance does not mean the feature is more important. In the house features example, number of bedrooms is as or more important than house area.

PCA on the other hand is a "dimension reduction" technique that is used when all entries of the original data vector are acquired using the same type of sensor / measuring tool. For example, image pixels. All CCD sensors are the same type in a camera. In such settings, it may be a valid assumption to claim that pixels with larger variance are more relevant features. When this assumption is true, then PCA can be used as a feature selection technique.

However PCA is not very useful if it is applied after feature normalization. Reason: feature norm is ensuring the new features all have unit variance. It can still be useful to de-correlate.

## 12.4    Practical information: How to implement PCA

### 12.4.1    Reduced dimension, $r$, is specified

**INPUT: $m \times n$ feature vectors' matrix $X$, chosen rank $r$.**

1. Compute $\hat{\mu} = \frac{1}{m} \sum_i x^{(i)}$

2. For $i = 1, 2, \ldots, m$, let $z^{(i)} = x^{(i)} - \hat{\mu}$ and let

$$Z = \begin{bmatrix} -(z^{(1)})^T- \\ -(z^{(2)})^T- \\ \vdots \\ -(z^{(m)})^T- \end{bmatrix} \tag{5}$$

   OR: compute directly

$$Z = X - \mathbf{1}_m(\mathbf{1}_m^T X / m)$$

   where $\mathbf{1}_m = [1\ 1\ \ldots\ 1]^T$ is an $m$-length vector of ones.

3. Compute the singular value decomposition (SVD) of $Z$ and set $V$ equal to the $r$ right singular vectors with the largest singular values (called "top $r$" right singular vectors), i.e., compute

$$Z \stackrel{SVD}{=} U_{full} S_{full} V_{full}^T$$

   where $U_{full}, V_{full}$ are unitary matrices and $S_{full}$ is an $m \times n$ diagonal matrix with non-negative entries (singular values) arranged in decreasing order of magnitude.

   Set $V = V_{full}(:, 1:r)$ in MATLAB notation (set $V$ equal to first $r$ columns of $V_{full}$).

   Thus $V$ is the $n \times r$ matrix whose columns span the computed *principal subspace*.

   - We can also obtain $V$ as the eigenvectors with the $r$ largest eigenvalues (called top $r$ eigenvectors) of $Z^T Z$.

4. Project the original data $X$ to $range(V)$: compute $b^{(i)} = V^T x^{(i)}$ for each $i$ or equivalently, compute

$$B = XV$$

   Thus $B$ is $m \times r$. These are the new feature vectors in the reduced dimensional space.

**OUTPUT: $m \times r$ matrix $B$: reduced dim feature vectors; $n \times r$ matrix $V$: principal subspace.**

   - If the rank $r$ approximation of $X$ is needed, this is obtained as the $m \times n$ matrix

$$L = BV^T = XVV^T$$

### 12.4.2 Deciding $r$

So far we have assumed that the desired lower-dimension $r$ is given. However, there is no one correct way of deciding $r$ in practice. Two common approaches:

1. A common heuristic is the 99% or some other high-enough percent heuristic: retain all eigenvectors so the that variance in the reduced dimensional space is at least 99% of the total variance of the data. In other words, find the smallest value $r$ so that

$$\sum_{i=1}^{r} \sigma_i^2 \geq 0.99 \sum_{i=1}^{\min(m,n)} \sigma_i^2.$$

   There is nothing "special" about 99%, we could also use another percentage.

2. An alternative approach is to pick the $r$ that is best for the final application for which PCA is being used as a pre-processing step. We use simple or leave-one-out cross-validation to compute the value of $r$ that minimizes the test-MSE.

### 12.4.3 Details for the second approach

**This approach is most useful if the feature vectors are noisy. It is typically not useful if output is a noisy function of feature vec's; in that case the best $r$ will often be equal or very close to $n$. Of course, in practice, for a real dataset, this part is not clear (it's clear for simulated data).**

Consider any classification application. Given $\boldsymbol{X}$, $\boldsymbol{y}$ and $\boldsymbol{X}_{test}, \boldsymbol{y}_{test}$: obtain by splitting all available data into 10% testing and 90% training). OR BETTER: use leave-one-out cross-validation (most efficient but more difficult to explain here).

1. Compute $\boldsymbol{Z} = \boldsymbol{X} - (\boldsymbol{1}_m \boldsymbol{1}_m^T \boldsymbol{X}/m)$

2. Compute SVD of $\boldsymbol{Z}$
$$\boldsymbol{Z} \overset{SVD}{=} \boldsymbol{U}_{full} \boldsymbol{S}_{full} \boldsymbol{V}_{full}^T$$

   where $\boldsymbol{U}_{full}, \boldsymbol{V}_{full}$ are unitary matrices and $\boldsymbol{S}_{full}$ is an $m \times n$ diagonal matrix with non-negative entries (singular values) arranged in decreasing order of magnitude.

3. Loop over $r$ from 1 to $n$

   (a) Set $\boldsymbol{V} = \boldsymbol{V}_{full}(:, 1:r)$ in MATLAB notation (set $\boldsymbol{V}$ equal to first $r$ columns of $\boldsymbol{V}_{full}$).

   (b) Compute $\boldsymbol{B} = \boldsymbol{X}\boldsymbol{V}$

   (c) Use $\boldsymbol{B}$ where you would have used $\boldsymbol{X}$: if linear regression then use it to get $\hat{\theta} = (\tilde{\boldsymbol{B}}^T \tilde{\boldsymbol{B}})^{-1} \tilde{\boldsymbol{B}}^\top \boldsymbol{y}$, if logistic reg, then use it in the GD algorithm to find $\hat{\theta}$.

   (d) Compute test error

      i. Compute $\boldsymbol{B}_{test} = \boldsymbol{X}_{test} \boldsymbol{V}$

      ii. Use $\boldsymbol{B}_{test}$ where you would have used $\boldsymbol{X}_{test}$:
         - if linear regression, then obtain $\hat{\boldsymbol{y}}_{test} = \boldsymbol{B}_{test}\hat{\theta}$ and compute $NormalizedTestMSE(r) = \|\boldsymbol{y}_{test} - \hat{\boldsymbol{y}}_{test}\|_2^2 / \|\boldsymbol{y}_{test}\|_2^2$
         - if logistic regression, then again obtain $\hat{\boldsymbol{y}}_{test}$ as explained in the Logistic Regression section, followed by computing the misclassification percentage

4. Pick $r$ for which $NormalizedTestMSE(r)$ is the smallest.

### 12.4.4  Sample Py Code

Three ways to get the SVD. The last one directly gives you the top $k = 2$ singular vectors. In case of the others, some post processing is needed to get the top $r$ singular vectors

```python
import numpy as np
from scipy import linalg
from scipy.sparse import linalg as slinalg

x = np.array([[1,1,0,0,0],[0,0,1,1,0],[1,1,1,1,1]],dtype=np.float64)

npsvd = np.linalg.svd(x)
spsvd = linalg.svd(x)
sptop = slinalg.svds(x,k=2)
```

## 12.5  Theory – computational - 1

PCA solves

$$\min_{\boldsymbol{L}} \|\boldsymbol{X} - \boldsymbol{L}\|_F$$

Set $\boldsymbol{L} = \boldsymbol{B}\boldsymbol{V}^\top$ and simplify, then this becomes

$$\min_{\boldsymbol{B},\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} \|\boldsymbol{X} - \boldsymbol{B}\boldsymbol{V}^\top\|_F^2 = \min_{\boldsymbol{B},\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} \sum_i \|\boldsymbol{x}^{(i)} - \boldsymbol{V}\boldsymbol{b}^{(i)}\|^2$$

Now we can solve for $\boldsymbol{b}^{(i)}$ in closed form: $\boldsymbol{b}^{(i)} = \boldsymbol{V}^\top\boldsymbol{x}^{(i)}$ and substitute it in. Thus we get

$$\min_{\boldsymbol{B},\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} \|\boldsymbol{X} - \boldsymbol{B}\boldsymbol{V}^\top\|_F^2 = \min_{\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} \sum_i \|(\boldsymbol{I} - \boldsymbol{V}\boldsymbol{V}^\top)\boldsymbol{x}^{(i)}\|^2$$

Since

$$\|(\boldsymbol{I} - \boldsymbol{V}\boldsymbol{V}^\top)\boldsymbol{x}\|^2 = \|\boldsymbol{x}\|^2 + \|\boldsymbol{V}^\top\boldsymbol{x}\|^2 - 2\|\boldsymbol{V}^\top\boldsymbol{x}\|^2 = \|\boldsymbol{x}\|^2 - \|\boldsymbol{V}^\top\boldsymbol{x}\|^2$$

the min simplifies to

$$\max_{\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} \sum_i \|\boldsymbol{V}^\top\boldsymbol{x}^{(i)}\|^2$$

Since

$$\sum_i \|\boldsymbol{V}^\top\boldsymbol{x}^{(i)}\|^2 = \sum_i (\boldsymbol{x}^{(i)\top}\boldsymbol{V}\boldsymbol{V}^\top\boldsymbol{x}^{(i)}) = trace(\sum_i \boldsymbol{x}^{(i)\top}\boldsymbol{V}\boldsymbol{V}^\top\boldsymbol{x}^{(i)}) = trace(\boldsymbol{V}^\top(\sum_i \boldsymbol{x}^{(i)}\boldsymbol{x}^{(i)\top})\boldsymbol{V})$$

thus this further simplifies to

$$\max_{\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} trace(\boldsymbol{V}^\top(\sum_i \boldsymbol{x}^{(i)}\boldsymbol{x}^{(i)\top})\boldsymbol{V})$$

Let EVD of $\hat{\Sigma} = (\sum_i \boldsymbol{x}^{(i)}\boldsymbol{x}^{(i)\top})/m = \boldsymbol{V}\Lambda\boldsymbol{V}^\top$ with $\boldsymbol{U}$ orthonormal Suppose $r = 2$. This then simplifies to

$$\max_{\boldsymbol{v}_2\perp\boldsymbol{v}_1,:\|\boldsymbol{v}_2\|=1}[\max_{\boldsymbol{v}_1:\|\boldsymbol{v}_1\|=1} \boldsymbol{v}_1^\top(\sum_i \boldsymbol{x}^{(i)}\boldsymbol{x}^{(i)\top})\boldsymbol{v}_1 + \boldsymbol{v}_2^\top(\sum_i \boldsymbol{x}^{(i)}\boldsymbol{x}^{(i)\top})\boldsymbol{v}_2]$$

The first term is maximized by the top eigenvector, $\boldsymbol{v}_1$ of $\hat{\Sigma} = (\sum_i \boldsymbol{x}^{(i)}\boldsymbol{x}^{(i)\top})/m$ and its maximum value equals the top eigenvalue, $\lambda_1$. Thus,

$$\max_{\boldsymbol{V}:\boldsymbol{V}^\top\boldsymbol{V}=\boldsymbol{I}} trace(\boldsymbol{V}^\top\hat{\Sigma}\boldsymbol{V}) = \lambda_1 + \max_{\boldsymbol{v}_2\perp\boldsymbol{v}_1,:\|\boldsymbol{v}_2\|=1} \boldsymbol{v}_2^\top\boldsymbol{U}\Lambda\boldsymbol{U}^\top\boldsymbol{v}_2$$

Thus, this is solved by $\boldsymbol{v}_2 =$ the second eigenvector of $\hat{\Sigma}$. Proceeding similarly

$$\hat{\boldsymbol{V}} = \boldsymbol{V}_{1:r}$$

solves the maximization.

Eigenvectors of $\hat{\Sigma}$ are right singular vectors of $\boldsymbol{X}$.

Then setting $\hat{\boldsymbol{B}} = \boldsymbol{X}\hat{\boldsymbol{V}}^\top$ proves that

$$\hat{\boldsymbol{L}} = \boldsymbol{X}\boldsymbol{V}_{1:r}\boldsymbol{V}_{1:r}^\top = \boldsymbol{U}_{1:r}\boldsymbol{S}_{r\times r}\boldsymbol{V}_{1:r}^\top$$

($r$-SVD of $\boldsymbol{X}$) solves

$$\min_{\boldsymbol{L}} \|\boldsymbol{X} - \boldsymbol{L}\|_F$$

## 12.6   Theory – computational 2

PCA also solves

$$\min_{\boldsymbol{L} \text{ rank } r} \|\boldsymbol{X} - \boldsymbol{L}\|_2^2 = \min_{\boldsymbol{L} \text{ rank } r} \lambda_{\max}((\boldsymbol{X} - \boldsymbol{L})^T(\boldsymbol{X} - \boldsymbol{L}))$$

By Weyl-type inequality for singular values, $\sigma_i(M) \leq \sigma_i(M_2) + \|M - M_2\|$ for any singular value $i$ and any two matrices $M, M_2$. Thus,

$$\|\boldsymbol{X} - \boldsymbol{L}\|_2 \geq \sigma_i(\boldsymbol{X}) - \sigma_i(\boldsymbol{L})$$

Since $\boldsymbol{L}$ is rank $r$, it has only $r$ nonzero singular values. Thus $\sigma_{r+1}(\boldsymbol{L}) = 0$. And so

$$\|\boldsymbol{X} - \boldsymbol{L}\|_2 \geq \sigma_{r+1}(\boldsymbol{X}).$$

Since this is true for any rank $r$ matrix $\boldsymbol{L}$, it is also true for the minimizer, i.e.

$$\min_{\boldsymbol{L}:\text{rank } r} \|\boldsymbol{X} - \boldsymbol{L}\|_2 \geq \sigma_{r+1}(\boldsymbol{X})$$

Now if we can find a specific matrix $\boldsymbol{L}$ for which $\|\boldsymbol{X} - \boldsymbol{L}\|_2 = \sigma_{r+1}(\boldsymbol{X})$ that will be the minimizer (since the minimum value cannot be any smaller than this).

Set

$$\hat{\boldsymbol{L}} = \sum_{i=1}^r \sigma_i u_i v_i^T$$

($r$-SVD of $\boldsymbol{X}$), then $\boldsymbol{X} - \hat{\boldsymbol{L}} = \sum_{i=r+1}^{\min(m,n)} \sigma_i u_i v_i^T$ and so

$$\|\boldsymbol{X} - \hat{\boldsymbol{L}}\|_2 = \sigma_{r+1}$$

here $\sigma_i = \sigma_i(\boldsymbol{X})$.

## 12.7   Statistical optimality

First assume everything is zero mean.

1. PCA finds the subspace $\boldsymbol{V}$ ($\boldsymbol{V}$ is a matrix with orthonormal columsn that define the subspace) and the projected random vector $\boldsymbol{b}$ so that the expected value of the squared 2-norm of the reconstruction error $\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b}$ is minimized. Thus, it solves

$$\min_{\boldsymbol{b}\in\Re^r, \boldsymbol{V}\in\Re^{n\times r}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b}\|_2^2]$$

If we minimize over $\boldsymbol{b}$ first as a function of $\boldsymbol{V}$, then this is a standard Least Squares problem whose solution is

$$\boldsymbol{b} = (\boldsymbol{V}^T\boldsymbol{V})^{-1}\boldsymbol{V}^T\boldsymbol{x} = \boldsymbol{V}^T\boldsymbol{x} \quad \text{since } \boldsymbol{V}^T\boldsymbol{V} = \boldsymbol{I} \text{ in this case}$$

Thus, we need to solve

$$\min_{\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{V}\boldsymbol{V}^T\boldsymbol{x}\|_2^2]$$

Since $\boldsymbol{V}^T\boldsymbol{V} = \boldsymbol{I}$, thus $\|\boldsymbol{V}\boldsymbol{V}^T\boldsymbol{x}\|_2^2 = \|\boldsymbol{V}^T\boldsymbol{x}\|_2^2$ and so $\|\boldsymbol{x} - \boldsymbol{V}\boldsymbol{V}^T\boldsymbol{x}\|_2^2 = \|\boldsymbol{x}\|^2 - \|\boldsymbol{V}^T\boldsymbol{x}\|^2$. With this, the above is also equivalent to

$$\max_{\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \mathbb{E}[\|\boldsymbol{V}^T\boldsymbol{x}\|^2]$$

Using property of trace, this is further equivalent to

$$\max_{\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} trace(\boldsymbol{V}^T\mathbb{E}[\boldsymbol{x}\boldsymbol{x}^T]\boldsymbol{V})$$

This is another commonly use definition for PCA: PCA finds the $r$ directions of largest variance of the data. Here $\mathbb{E}[\boldsymbol{x}\boldsymbol{x}^T]$ is the covariance matrix.

2. PCA can also be understood as minimizing the worst-case (largest) expected reconstruction error in any direction: for direction $\boldsymbol{w}$, this is $|\boldsymbol{w}^T(\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b})|$

$$\min_{\boldsymbol{b},\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \max_{\boldsymbol{w}:\|\boldsymbol{w}\|_2=1} \mathbb{E}[(\boldsymbol{w}^T(\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b}))^2]$$

Using the properties of trace, this is equivalent to

$$\min_{\boldsymbol{b},\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \max_{\boldsymbol{w}:\|\boldsymbol{w}\|_2=1} \boldsymbol{w}^T\mathbb{E}[(\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b})(\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b})^T]\boldsymbol{w} = \min_{\boldsymbol{b},\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \lambda_{\max}(\mathbb{E}[(\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b})(\boldsymbol{x} - \boldsymbol{V}\boldsymbol{b})^T])$$

The last equality follows using the variational definition of the maximum eigenvalue.

In the nonzero mean case, we do either of the above for $(\boldsymbol{x} - \mu)$.

## 12.8   PCA de-correlates the data: what does it mean **

This claim depends on how the principal subspace is defined. In (statistical) theory, we are finding $\boldsymbol{V}$ and $\boldsymbol{b}$ that solves $\min_{\boldsymbol{b},\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \mathbb{E}[\|(\boldsymbol{x} - \mu) - \boldsymbol{V}\boldsymbol{b} + \boldsymbol{V}\boldsymbol{V}^T\mu)\|_2^2]$, i.e., it minimizes the expected value of the reconstruction error. Since the minimizer over $\boldsymbol{b}$ is $\boldsymbol{b} = \boldsymbol{V}^T\boldsymbol{x}$, we are actually finding $\boldsymbol{V}$ that solves $\min_{\boldsymbol{V}:\boldsymbol{V}^T\boldsymbol{V}=\boldsymbol{I}} \mathbb{E}[\|(\boldsymbol{x} - \mu) - \boldsymbol{V}\boldsymbol{V}^T(\boldsymbol{x} - \mu)\|_2^2]$ and then setting $\boldsymbol{b} = \boldsymbol{V}^T\boldsymbol{x}$. If we can find this $\boldsymbol{V}$, then the computed lower dimensional r.v. $\boldsymbol{b}$ satisfies

$$\mathbb{E}[(\boldsymbol{b} - \mathbb{E}[\boldsymbol{b}])(\boldsymbol{b} - \mathbb{E}[\boldsymbol{b}])^T] \text{ is diagonal}$$

i.e.

$$\mathbb{E}[(\boldsymbol{b} - \mathbb{E}[\boldsymbol{b}])_j(\boldsymbol{b} - \mathbb{E}[\boldsymbol{b}])_k] = 0$$

for $j \neq k$.

*Proof: assume everything is zero mean for ease of writing. We have* $\boldsymbol{b} = \boldsymbol{V}'\boldsymbol{x}$ *so* $b_j = \boldsymbol{v}_j'\boldsymbol{x}$, *thus,* $\mathbb{E}[b_j b_k] = \boldsymbol{v}_j'\mathbb{E}[\boldsymbol{x}\boldsymbol{x}']\boldsymbol{v}_k$. *As explained in previous section,* $\boldsymbol{V}$ *is the matrix of top* $r$ *eigenvectors of* $\mathbb{E}[\boldsymbol{x}\boldsymbol{x}']$, *i.e., that* $\mathbb{E}[\boldsymbol{x}\boldsymbol{x}'] = \boldsymbol{V}\Sigma\boldsymbol{V}' + \boldsymbol{V}_\perp\Sigma_\perp\boldsymbol{V}_\perp'$. *Hence,* $\mathbb{E}[b_j b_k] = \boldsymbol{v}_j'(\boldsymbol{V}\Sigma\boldsymbol{V}' + \boldsymbol{V}_\perp\Sigma_\perp\boldsymbol{V}_\perp')\boldsymbol{v}_k = \boldsymbol{v}_j'(\boldsymbol{V}\Sigma\boldsymbol{V}')\boldsymbol{v}_k = 0$.

In practice, we cannot find above but only its data-based (empirical) approximation. Hence in practice, we are always finding $\boldsymbol{V}$ as the top $r$ right singular vectors of $\boldsymbol{Z} = \boldsymbol{X} - \boldsymbol{1}_m\boldsymbol{1}_m^T\boldsymbol{X}/m$ (recall: where $\boldsymbol{1}_m = [1\ 1\ \ldots\ 1]^T$ is an $m$-length vector of ones). With this choice of $\boldsymbol{V}$, "uncorrelated" means the following:

$$\sum_{i=1}^m (\boldsymbol{b}^{(i)})_j(\boldsymbol{b}^{(i)})_k = 0$$

for $j \neq k$. In other words, the columns of the matrix $\boldsymbol{B}$ are mutually orthogonal.

*Proof: same basic idea as above.*

# 13   Unsupervised Learning: Clustering

Use ML-cs229- notes on Clustering to see the figures.

## 13.1   Problem

Given unlabeled data/features $\boldsymbol{x}^{(i)}$, $i = 1, 2, \ldots, m$, the goal is to partition the dataset into "cohesive" clusters (all points in same cluster are "close" while those in different clusters are "far"). Suppose we want to partition into $k$ clusters. Then the goal can also be stated as: for each $i$, find the class label $y^{(i)}$ (this can take values from $\{1, 2, \ldots, k\}$).

## 13.2   k-means clustering

The goal is to find the class labels $y^{(i)}$ for each data point and the cluster centers so that the following cost is minimized

$$J(\mu_j, j = 1, 2, \ldots, k, \boldsymbol{y}) = \sum_{i=1}^{m} \|\boldsymbol{x}^{(i)} - \mu_{y^{(i)}}\|_2^2$$

The original k-means clustering algorithm provides an Alternating-Minimization (AltMin) algorithm to minimize the above cost.

1. Initialize cluster centers $\hat{\mu}_1, \hat{\mu}_2, \ldots, \hat{\mu}_k$. Can do random init.

2. Repeat

   (a) For each $i = 1, 2, \ldots, m$, find the class labels

   $$\hat{y}^{(i)} = \arg \min_{j=1,2,\ldots,k} \|\boldsymbol{x}^{(i)} - \hat{\mu}_j\|_2^2.$$

   (b) Update cluster centers: for each $j = 1, 2, \ldots, k$, compute

   $$\hat{\mu}_j = \frac{\sum_{i=1}^{m} \mathbb{1}(\hat{y}^{(i)} == j)\boldsymbol{x}^{(i)}}{\sum_{i=1}^{m} \mathbb{1}(\hat{y}^{(i)} == j)}$$

   (the above is a solution to $\arg \min_{\mu_j, j=1,2,\ldots,k} \sum_{i=1}^{m} \|\boldsymbol{x}^{(i)} - \mu_{yhat^{(i)}}\|_2^2$)

   until cluster center estimates do not change much, i.e., until $\max_j(\|\hat{\mu}_j^{(t+1)} - \hat{\mu}_j^{(t)}\|/\|\hat{\mu}_j^{(t)}\|) < threshold$ where $threshold = 0.001$ or some small fraction.

3. IMPROVED VERSION: repeat above algorithm $N$ times with different random init's each time. For each repeat, compute the cost function value $J(\hat{\mu}_j, j = 1, 2, \ldots, k, \hat{y})$. Keep the output of the repeat with the smallest cost.

IMPROVEMENT 2: one can replace the regular Euclidean distance to the cluster center by any other distance that is more relevant to the application. As an example, if it is known that different features are likely to have significantly different variances, one could init with $\hat{\Sigma}_j = \boldsymbol{I}$, replace $\arg \min_{j=1,2,\ldots,k} \|\boldsymbol{x}^{(i)} - \hat{\mu}_j\|_2$ in step 2a by

$$\arg \min_{j=1,2,\ldots,k} (\boldsymbol{x}^{(i)} - \hat{\mu}_j)\hat{\Sigma}_j^{-1}(\boldsymbol{x}^{(i)} - \hat{\mu}_j)$$

and in step 2b, update

$$\hat{\Sigma}_j = \frac{\sum_{i=1}^{m} \mathbb{1}(\hat{y}^{(i)} == j)(\boldsymbol{x}^{(i)} - \hat{\mu}_j)(\boldsymbol{x}^{(i)} - \hat{\mu}_j)^T}{\sum_{i=1}^{m} \mathbb{1}(\hat{y}^{(i)} == j)}$$

The above is an example of AltMin. See Sec. 5.4.

### 13.3  Probabilistic Model (Generative Model) for Clustering: Gaussian Mixture Model (GMM)

The Gaussian Mixture Model or GMM is a common way to specify a clustering problem. $\boldsymbol{x}$ follows the GMM with $k$ components means that

$$p(\boldsymbol{x};\theta) = \sum_{j=1}^{k} \mathcal{N}(\boldsymbol{x};\mu_j,\Sigma_j)\phi_j$$

where $\phi_j$'s are the mixture weights (probability of $x$ coming from the $j$-th class in the mixture) and thus $\sum_{j=1}^{k}\phi_j = 1$. We often refer to $j$ as the class labels.

GMM: means that $\boldsymbol{x}$ is generated from class $j$ with probability $\phi_j$, and given that it is generated from class $j$, it follows a Gaussian distribution with mean $\mu_j$ and covariance $\Sigma_j$.

The model assumed by Gaussian discriminant analysis (GDA) covered earlier and in HW 2 is also GMM. See Sec 5.2. Except there we had labeled data (for each training data point, the class label was available), so it was easy to "learn" the model parameters. As a result, the learning of $\phi_j$'s was decoupled from the learning of the mean and covariances. As a result we in fact got closed form expressions for the MLE.

Here we do not have class labels and thus the learning problem is more difficult. We need to use Gradient Descent or some other iterative approach.

### 13.4  EM algorithm for MLE for Gaussian Mixture Model

A popular approach for MLE for the GMM is the EM algorithm. It uses AltMin to minimize a lower bound on the MLE cost. This is fully explained and derived in old version of these nodes ML-algorithms-full.pdf.

## 14  Reliability of ML estimates, Regularization / Feature Selection / Picking reduced dimension $r$ to minimize test-error / generalization error (optimize Bias-Variance Tradeoff)

### 14.1  Reliability of an output

Linear regression predicts a real-valued scalar where as everything we learnt after that predicts a class label (solves a classification problem). Given a query, we can always obtain a prediction. But the other important question to answer is : how reliable is the prediction we obtained. The answer to this question depends on

- the problem itself, e.g., in case of classification by GDA, if the two class means are very close, it is not easy to distinguish the classes. More precisely what matters is how close the class means along a given direction compared to the standard deviation along that direction. Practically this means the following: it easier to distinguish dog pictures from human pictures than from cat pictures

  Similarly for a regression problem, the amount of modeling error $e$ or its variance decides how good the prediction is.

- the number of training data points, and how well the training and test data match (this decides quality of learnt model). In most of what we learn, it is assumed that training and test data are genarated from the same distribution, but in real life this may not be true.

- the specific query: if the query image is of a fluffy cat that may look dog-like, then it is hard to reliably provide a correct classification.

- the last problem can be partly addressed by changing the learning algorithm (the assumed model on the data).

## 14.2 Bias, Variance, Model mismatch errors and Bias Variance Tradeoff

Consider a generative model: suppose that $y, \boldsymbol{x}$ satisfy

$$y = f(\boldsymbol{x}) + e, \ \mathbb{E}[e] = 0, \mathbb{E}[e^2] = \sigma^2,$$

with $e$ being zero mean "modeling error"/"noise" that is independent of $\boldsymbol{x}$, it is also independent for each data point $y_i$, and $\theta$ being the model parameters. We do not know $f(.)$. But we assume $f(\boldsymbol{x}) \approx h_\theta(\boldsymbol{x})$ i.e., we try to "model" it as $\hat{y} = \hat{f}(x) = h_{\hat{\theta}}(\boldsymbol{x})$, e.g., in linear regression, $\hat{f}(x) = h_{\hat{\theta}}(\boldsymbol{x}) = \hat{\theta}^T \boldsymbol{x}$ with $\hat{\theta}$ estimated by Maximum Likelihood estimation (MLE) as described earlier using training data. In logistic regression, $\hat{f}(x) = h_{\hat{\theta}}(\boldsymbol{x}) = g(\hat{\theta}^T \boldsymbol{x})$ with $g(.)$ being the sigmoid function and $e$ is randomness so that $p(y = 1) = h_\theta(\boldsymbol{x})$. We again estimate $\theta$ by MLE using training data. MLE for linear reg has a closed form solution. For Log Reg, it needs to use GD.

MLE uses "training data"

$$(y^{(i)}, \boldsymbol{x}^{(i)}), i = 1, 2, \ldots, m$$

The question is how good is my learnt model (in terms of mean squared error or MSE on test data or some other metric), i.e., for a test query $\boldsymbol{x}$, what is $\mathbb{E}[(y - \hat{f}(\boldsymbol{x}))^2]$ and what can we do to improve it? Let us assume MSE as the error metric.

We define Test-MSE as

$$\text{Test-MSE} := \mathbb{E}[(y - \hat{y})^2] = \mathbb{E}[(y - \hat{f}(\boldsymbol{x}))^2] = \mathbb{E}[(f(\boldsymbol{x}) + e - \hat{f}(\boldsymbol{x}))^2]$$

over test data, i.e., $\mathbb{E}[.] \equiv \mathbb{E}_{\text{test data}}[.]$. Usually $\boldsymbol{x}$ is treated as a constant, so then the expected value is over the distribution of the noise $e$.

Since $e$ is test-data noise, it is independent of $\hat{f}(\boldsymbol{x}) = h_{\hat{\theta}}(\boldsymbol{x})$ since $\hat{\theta}$ was estimated using training data. Also, by assumption, $e$ it is independent of $f(x)$. Thus, we have

$$\begin{aligned}
\text{Test-MSE} &= \mathbb{E}[e^2] + \mathbb{E}[(f(\boldsymbol{x}) - \hat{f}(\boldsymbol{x}))^2], \quad \hat{f}(\boldsymbol{x}) = h_{\hat{\theta}}(\boldsymbol{x}) \\
&= \mathbb{E}[e^2] + (\mathbb{E}[f(\boldsymbol{x}) - \hat{f}(\boldsymbol{x})])^2 + Variance[f(\boldsymbol{x}) - \hat{f}(\boldsymbol{x})] \\
&= \sigma^2 + Bias^2 + Variance
\end{aligned}$$

where $Variance(Z) := E[(Z - E[Z])^2]$. The first term, $\sigma^2$, depends on how noisy my data is. The second two terms depend on the "assumed model" and how well its parameters are estimated. Bias depends on how good my assumed model is, while variance depends on how well I can estimate $\theta$. Usually

$$Bias = f(\boldsymbol{x}) - h_\theta(\boldsymbol{x}), \ Variance = Var(h_\theta(\boldsymbol{x}) - h_{\hat{\theta}}(\boldsymbol{x}))$$

A model with more parameters "generally" will have lower bias. But for a fixed value of $m$ (training data size), variance in estimating $\theta$ will be higher.

## 14.3 Approximating Test Error / Generalization Error / Test-MSE in practice

While we can write things as above, it is *not* actually possible to compute the above decomposition for test data.

**Simple approach:** All one can do is the following: for a given model, one can approximate Test-MSE using the following approach

- Split available training data into training and test data: in other words do not use all $m$ data points to train, split them so that $m = m_{train} + m_{test}$.

- Use the $m_{train}$ data points to train, i.e. to estimate $\theta$ for the assumed model

- Approximate Test-MSE by $\frac{1}{m_{test}} \sum_{j=1}^{m_{test}} (y_j - h_{\hat{\theta}}(\boldsymbol{x}_j))^2$

The above is one way to do what is called "Cross-Validation". Typically, one uses $m_{test} = 0.25m$ and $m_{train} = 0.75m$ or similar.

**Leave-one-out Cross-Validation:** Do above with $m_{test} = 1$, and $m_{train} = m - 1$, but repeat the procedure $m$ times with a different test sample, and compute error each time. Average it to get Test-MSE. See Cross Validation section earlier.

## 14.4 How to reduce Test Error a.k.a. Regularization a.k.a. Feature Selection: reduce variance, not increase bias too much

Goal: Try to reduce variance, while not increasing the bias too much.

With each new intervention, compute the Test-MSE as explained above, see if it gets reduced or not. In general as you keep reducing the number of features that are being used, variance reduces but bias increases. Initially when you begin reducing the number, the reduction in variance is significant but the increase in bias is negligible. So Test MSE will decrease. This will happen until the bias starts increasing (Test MSE starts increasing). Use the number of features that results in the smallest Test MSE.

1. Regularization on $\theta$:

    (a) Stepwise regression; start with no features; at each iteration, select the one feature that is the "best" to add from remaining set. Define "best" using some heuristic (say t-test etc) or by model fitting using each new feature.

    (b) Use domain knowledge such as assuming entries of $\theta$ are in decreasing order of magnitude. If this is known, then you just use the first $r$ features. Vary $r$ to pick which option reduces test data error the most.

    (c) Sparsity: Model $\theta$ as being sparse (only a few features are important but we do not know which ones), this is a generalization of the assumption used in item 1. In this case, we add $\|\theta\|_1$ into the cost function for learning $\theta$.

    In general, if the cost function is $J(\theta)$, then this approach says we replace it by $J(\theta) + \lambda\|\theta\|_1$

    In case of linear regression this becomes "Lasso"

    $$\min_{\theta} \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2 + \lambda\|\theta\|_1$$

    **Do not try to write your own GD code for this. Use cvxopt or cvxpy etc.**
    **Vary $\lambda$ as multiples of $m$, so $\lambda = \beta m$, with $\beta = 0.01, 0.1, 1, 10$ etc**
    More information in Sec. 6.

    (d) $\ell_2$ norm: Use the assumption that all features are important and no feature is too important (no weight can be too large): then use the $\ell_2$ regularization, i.e. add $\|\theta\|^2$ into the cost function for learning $\theta$. In general, if the cost function is $J(\theta)$, then this approach says we replace it by $J(\theta) + \lambda\|\theta\|^2$

    In case of linear regression this becomes "Ridge Regression"

    $$\min_{\theta} \|\boldsymbol{y} - \boldsymbol{X}\theta\|^2 + \lambda\|\theta\|^2$$

    **This cost has closed form solution, given next.**

    $$\hat{\theta} = (\boldsymbol{X}^\top\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}^\top\boldsymbol{y}$$

    More information in Sec. 6.

    (e) Prior estimate $\theta_0$ of $\theta$ given: Suppose prior knowledge is available that $\theta$ is close to a given vector $\theta_0$, then add the $\|\theta - \theta_0\|^2$ into the cost function for finding $\theta$. See Sec. 6.4 under Use of prior knowledge and Recursive LS.

41

2. Regularization on input features' matrix $\boldsymbol{X}$ using PCA or Feature normalization: PCA is a dimension reduction technique. It is used for feature selection under the assumption that "features having the most variation are the most important ones". Also of course PCA helps de-correlate the features. A typical setting where this is useful is when the raw feature vector is all image pixels. There are too many pixels and pixels are highly correlated. As explained earlier, PCA is a good pre-processing approach to regularization/feature-selection before doing classification (GDA or Log Reg or SVMs). *PCA does not do much if features are first "normalized" and then PCA is applied, see Sec. 12.3.1 and 12.3.2.*

3. In case of Clustering, reducing $k$ (number of classes) will reduce the variance, increasing it will reduce the bias.

4. Naive Bayes assumption described earlier is another way to reduce the number of parameters. This is used in generative or Bayesian learning settings where the feature vector itself is modeled probabilistically.

With each new intervention, compute the Test-MSE as explained above, see if it gets reduced or not. In general as you keep reducing the number of features that are being used, variance reduces but bias increases. Initially when you begin reducing the number, the reduction in variance is significant but the increase in bias is negligible. So Test MSE will decrease. This will happen until the bias starts increasing (Test MSE starts increasing). Use the number of features that results in the smallest Test MSE.

# A    Optimization overview

Consider differentiable cost functions.

   Unconstrained opt: solve

$$\min_{\theta} f(\theta)$$

A necessary condition for $\hat{\theta}$ to be a minimizer is

$$\frac{\partial f}{\partial \theta}(\hat{\theta}) = \mathbf{0}$$

   Constrained opt: solve

$$\min_{\theta} f(\theta), \;\; s.t. \; \boldsymbol{h}(\theta) = \mathbf{0}, \; \boldsymbol{g}(\theta) \preceq 0$$

Here $\boldsymbol{z} \preceq \mathbf{0}$ means $\boldsymbol{z}_i \leq 0$ for all entries of $i$ of the vector $\boldsymbol{z}$.

   Under certain "regularity conditions" (sometimes called constraint qualifications), one set of necessary conditions for $\hat{\theta}$ to be a local minimizer is the *KKT conditions*.

   Let $\boldsymbol{\mu}, \boldsymbol{\lambda}$ be Lagrange multipliers for the inequality and equality constraints respectively. Define the Lagrangian as

$$L(\theta, \boldsymbol{\lambda}, \boldsymbol{\mu}) := f(\theta) + \boldsymbol{\lambda}^{\top} \boldsymbol{h}(\theta) + \boldsymbol{\mu}^{\top} \boldsymbol{g}(\theta)$$

   Necessary conditions for $\hat{\theta}$ to be a minimizer:

1. Derivative of Lagrangian is zero

$$\frac{\partial L}{\partial \theta}(\hat{\theta}) = \mathbf{0} \;\;\Leftrightarrow\;\; \frac{\partial f}{\partial \theta}(\hat{\theta}) + \boldsymbol{\lambda}^{\top} \frac{\partial \boldsymbol{h}}{\partial \theta}(\hat{\theta}) + \boldsymbol{\mu}^{\top} \frac{\partial \boldsymbol{g}}{\partial \theta}(\hat{\theta}) = \mathbf{0}$$

2. The problem constraints hold – called Primal feasibility

$$\boldsymbol{h}(\hat{\theta}) = \mathbf{0}, \;\; \text{and} \;\; \boldsymbol{g}(\hat{\theta}) \preceq 0$$

3. The Lagrange multipliers for the inequality constraints are non-negative – called Primal feasibility

$$\boldsymbol{\mu} \succeq 0$$

4. Complementary slackness holds: if an inequality constraint is *not* satisfied with equality, then the corresponding Lag multiplier is zero.

$$\boldsymbol{\mu}_i \boldsymbol{g}_i(\hat{\theta}) = 0 \;\; \text{for all } i = 1, 2, \ldots, m$$

   In some texts this condition is written in a more compact (but more confusing-looking) fashion as

$$\boldsymbol{\mu}^{\top} \boldsymbol{g}(\theta) = \mathbf{0}$$

   This and previous one are equivalent because we are assuming $\boldsymbol{g}(\theta) \preceq 0$ and $\boldsymbol{\mu} \succeq 0$. This implies each term in the summation $\boldsymbol{\mu}^{\top} \boldsymbol{g}(\theta)$ is $\leq 0$. The only way the sum becomes zero is if each term is zero.

## A.1    Convex opt problem

An opt problem is convex if the cost is convex and the feasibility set (set of $\theta$s in the constraint set) is convex. One can show that this is the same as requiring

1. the inequality constraint functions $\boldsymbol{g}(\theta)$ are convex functions, i.e. $\boldsymbol{g}_i(\theta)$ is convex function for each $i$

2. the cost function $f(\theta)$ is convex function

3. equality constraints $\boldsymbol{h}(\theta)$ are affine, i.e., $\boldsymbol{h}(\theta) = \boldsymbol{A}\theta + \boldsymbol{b}$ for some matrix $\boldsymbol{A}$ and vector $\boldsymbol{b}$.

### A.1.1 Convex function

A function is convex if for any $0 \leq \alpha \leq 1$

$$f(\alpha \theta_1 + (1 - \alpha)\theta_2) \leq \alpha f(\theta_1) + (1 - \alpha)f(\theta_2)$$

### A.1.2 Convex set

A set is convex if $\theta_1, \theta_2$ in the set implies that $\alpha \theta_1 + (1 - \alpha)\theta_2$ is also in the set for any $0 \leq \alpha \leq 1$.

## A.2 Matrix calculus: key points

1. $\nabla_{\boldsymbol{x}} \boldsymbol{b}^\top \boldsymbol{x} = \boldsymbol{b}$

2. For $\boldsymbol{A}$ symmetric, $\nabla_{\boldsymbol{x}} (\boldsymbol{x}^\top \boldsymbol{A} \boldsymbol{x}) = 2\boldsymbol{A}\boldsymbol{x}$

3. $\nabla_{\boldsymbol{A}} \log det(\boldsymbol{A}) = \boldsymbol{A}^{-1}$

4. $\nabla_{\boldsymbol{B}} tr(\boldsymbol{A}\boldsymbol{B}) = \boldsymbol{A}^\top$

### A.2.1 Application of the above to derive MLE estimates for $K$ class GDA

Goal: find

$$\arg \max_{\phi_k, \mu_k, k=1,..K, \Sigma \text{ s.t. } \sum_k \phi_k = 1} \prod_{i=1}^{m} p(\boldsymbol{y}^{(i)})p(\boldsymbol{x}^{(i)}|\boldsymbol{y}^{(i)})$$

with

$$p(y) := \prod_{k=1}^{K} \phi_k^{\mathbb{1}(y=k)}, \quad p(\boldsymbol{x}|y) = \prod_{k=1}^{K} \mathcal{N}(\boldsymbol{x}; \mu_k, \Sigma)^{\mathbb{1}(y=k)}$$

Let

$$\theta = \{\phi_k, \mu_k, k = 1, ..K, \Sigma\}$$

This is equivalent to finding

$$\arg \min_{\theta} J(\theta) \text{s.t.} \sum_k \phi_k = 1 \tag{6}$$

$$J(\theta) := -\log[\prod_{i=1}^{m} p(\boldsymbol{y}^{(i)}; \theta)p(\boldsymbol{x}^{(i)}|\boldsymbol{y}^{(i)}; \theta)] \tag{7}$$

Define the Lagrangian as

$$L(\theta, \lambda) = -\sum_i \log[p(\boldsymbol{y}^{(i)}; \theta)p(\boldsymbol{x}^{(i)}|\boldsymbol{y}^{(i)}; \theta)] + \lambda(\sum_{k=1}^{K} \phi_k - 1)$$

$$L(\theta, \lambda) = -\sum_i \log \left[ \prod_{i=1}^{m} \prod_{k=1}^{K} \phi_k^{\mathbb{1}(y^{(i)}=k)} \prod_{i=1}^{m} \prod_{k=1}^{K} \mathcal{N}(\boldsymbol{x}^{(i)}; \mu_k, \Sigma)^{\mathbb{1}(y^{(i)}=k)} \right] + \lambda(\sum_{k=1}^{K} \phi_k - 1)$$

This simplifies to $L(\theta, \lambda) =$

$$-\sum_i \sum_k \mathbb{1}(\boldsymbol{y}^{(i)} = k) \log \phi_k + \sum_i \sum_k \mathbb{1}(\boldsymbol{y}^{(i)} = k)[(\boldsymbol{x}^{(i)} - \mu_k)^\top \Sigma^{-1}(\boldsymbol{x}^{(i)} - \mu_k) + \log det(\Sigma)] + C + \lambda(\sum_{k=1}^{K} \phi_k - 1)$$

Using the necessary conditions given earlier (KKT conditions), and key points from matrix calculus, the minimizer needs to satisfy

$$\nabla_\theta L(\theta, \lambda) = 0, \quad \text{and} \quad \sum_{k=1}^{K} \phi_k = 1$$

This simplifies to

$$-\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k)\frac{1}{\phi_k} + \lambda = 0, \quad \text{and} \quad \sum_{k=1}^{K} \phi_k = 1$$

and

$$2\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k)\Sigma^{-1}(\boldsymbol{x}^{(i)} - \mu_k) = \boldsymbol{0}$$

and using $(\boldsymbol{x}^{(i)} - \mu_k)\Sigma^{-1}(\boldsymbol{x}^{(i)} - \mu_k) = trace((\boldsymbol{x}^{(i)} - \mu_k)^\top \Sigma^{-1}(\boldsymbol{x}^{(i)} - \mu_k)) = trace((\boldsymbol{x}^{(i)} - \mu_k))(\boldsymbol{x}^{(i)} - \mu_k))^\top \Sigma^{-1})$,

$$\sum_i \sum_k \mathbb{1}(\boldsymbol{y}^{(i)} = k)[(\boldsymbol{x}^{(i)} - \mu_k)(\boldsymbol{x}^{(i)} - \mu_k)^\top - \Sigma] = 0$$

Simplifying the last three equations gives the final solution:

$$\phi_k = \frac{1}{\lambda}\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k) = \frac{\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k)}{\sum_i 1} = \frac{\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k)}{m}$$

$$\mu_k = \frac{1}{\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k)}\sum_i \mathbb{1}(\boldsymbol{y}^{(i)} = k)\boldsymbol{x}^{(i)}$$

$$\Sigma = \frac{1}{m}\sum_i \sum_k \mathbb{1}(\boldsymbol{y}^{(i)} = k)[(\boldsymbol{x}^{(i)} - \mu_k)(\boldsymbol{x}^{(i)} - \mu_k)^\top$$

# B    Multi-class Log Reg

## B.1    Multi-class Logistic Regression**

Multi-class logistic regression (LR) is known by a variety of other names, including polytomous LR, Multinomial LR, softmax regression, multinomial logit (mlogit), the maximum entropy (MaxEnt) classifier, and the conditional maximum entropy model.

We now have $K > 2$ classes. Recall that, for $K = 2$, the model is

$$\Pr(y = 1) = \frac{1}{1 + \exp(-\theta^\top \boldsymbol{x})} = \frac{\exp(\theta^\top \boldsymbol{x})}{1 + \exp(\theta^\top \boldsymbol{x})}$$

and

$$\Pr(y = 0) = 1 - \frac{\exp(\theta^\top \boldsymbol{x})}{1 + \exp(\theta^\top \boldsymbol{x})}$$

To extend this to $K$ classes labeled $k = 0, 1, \ldots, K - 1$, instead of a single "regression vector" $\theta$, we now use $K - 1$ "regression" vectors $\beta_k, k = 1, 2, \ldots K - 1$ and use the model:

$$\Pr(y = k) = \frac{\exp(\beta_k^\top \boldsymbol{x})}{1 + \sum_{k'=1}^{K-1}\exp(\beta_{k'}^\top \boldsymbol{x})}, \quad k = 1, 2, \ldots, K - 1$$

and

$$\Pr(y = 0) = \frac{1}{1 + \sum_{k'=1}^{K-1}\exp(\beta_{k'}^\top \boldsymbol{x})}$$

If we let $\beta_0 = \mathbf{0}$ the $y = 0$ case can be interpreted as

$$\Pr(y = 0) = \frac{\exp(\beta_0^\top \boldsymbol{x})}{1 + \sum_{k'=1}^{K-1} \exp(\beta_{k'}^\top \boldsymbol{x})}, \text{ with } \beta_0 = \mathbf{0}$$

With setting $\beta_0 = \mathbf{0}$, observe that $1 = \exp(\beta_0^\top \boldsymbol{x})$. Replace this in both numerator and denominator of above expressions. Then, the above can be simplified to

$$\Pr(y = k) = \frac{\exp(\beta_k^\top \boldsymbol{x})}{\sum_{k'=0}^{K-1} \exp(\beta_{k'}^\top \boldsymbol{x})}, \ k = 1, 2, \ldots, K-1$$

- Learning the parameters:

  - Need to estimate $\theta := \{\beta_1, \beta_2, \ldots, \beta_{K-1}\}$.
  - Solution 1: define the joint log-likelihood function over all $m$ training data points and maximize it (minimize its negative).
  - This solution often does not work too well (too many parameters and one may start becoming too large).
  - Solution 2: add regularization on the $\beta_k$'s, e.g. add $\lambda \sum_{k=1}^{K-1} \|\beta_k\|^2$ to the cost function and then minimize it
  - Basic GD may not work easily in this complicated setting. To explore this case, see what Python or MATLAB's built in toolboxes do for your data.

- Classification: find the $k$ for which $\Pr(y = k)$ is highest.
  Equivalently, taking log of the probabilities, using $\beta_0 = \mathbf{0}$ (zero vector), and ignoring the denominator, this translates to

$$\hat{y} = \arg \max_{k=0,\ldots,K-1} \beta_k^\top \boldsymbol{x}$$

## B.2   Another way to understand multi-logit

Suppose there are $K$ classes. Instead of letting $y$ be a scalar that takes $K$ possible values, we instead let $\boldsymbol{y}$ be a $K$-length "one hot" vector. This means $\boldsymbol{y} = \boldsymbol{e}_k$ if the class label is $k$. Here $\boldsymbol{e}_k$ is a vector with all zeros except at the $k$-th index.

Also we define the softmax function as one that takes as input a $K$-length real vector with any values and outputs a $K$-length vector with each entry being between zero and one (being a probability), as follows:

$$\boldsymbol{y} = \text{softmax}(\boldsymbol{z}) \text{ with } \ \boldsymbol{y}_k = \frac{\exp(\boldsymbol{z}_k)}{\sum_{k'=1}^{K} \exp(\boldsymbol{z}_{k'})}, k = 1, 2, \ldots, K$$

Then, the multi-logit model can be described as follows. Let

$$\boldsymbol{B} := [\beta_1, \beta_2, \ldots, \beta_K] \text{ this is a } n \times K \text{ matrix}$$

with one of them say $\beta_1 = \mathbf{0}$. Also let

$$\boldsymbol{p}_k(\boldsymbol{x}; \boldsymbol{B}) := \Pr(\boldsymbol{y} = \boldsymbol{e}_k; \boldsymbol{x}, \boldsymbol{B})$$

Then,

$$\boldsymbol{p}(\boldsymbol{x}; \boldsymbol{B}) = \text{softmax}(\boldsymbol{B}^\top \boldsymbol{x})$$

Given training data, the max likelihood estimate (MLE) in this case can be obtained by minimizing the following negative log likelihood function

$$J(\boldsymbol{B}) := -\sum_{i=1}^{m} \sum_{k=1}^{K} \boldsymbol{y}_k^{(i)} \log \text{softmax}(\boldsymbol{B}^\top \boldsymbol{x}^{(i)})_k$$

The above is can be interpreted in terms of cross-entropy loss between two probability distributions: for two discrete probability vectors $\boldsymbol{q}, \boldsymbol{p}$,

$$CE(\boldsymbol{q}, \boldsymbol{p}) = -\sum_k \boldsymbol{q}_k \log \boldsymbol{p}_k$$

Thus,

$$J(\boldsymbol{B}) = \sum_{i=1}^{m} CE(\boldsymbol{y}^{(i)}, \boldsymbol{p}(\boldsymbol{x}^{(i)}; \boldsymbol{B}))$$

Notice the vector $\boldsymbol{y}$ can be interpreted as a probability.

This cost function $J(\boldsymbol{B})$ is not convex and complicated. If we minimize it as is, we may end up with very bad minimizers such as all but one $\beta_k$ being all zeros. In such cases, there is a need to regularize. Most common and default approach in Python's sci-kit-learn is Ridge Regularization: to add $\sum_k ||\beta_k||^2/C = ||\boldsymbol{B}||_F^2/C$. Minimize

$$J'(\boldsymbol{B}) = \sum_{i=1}^{m} CE(\boldsymbol{y}^{(i)}, \boldsymbol{p}(\boldsymbol{x}^{(i)}; \boldsymbol{B})) + ||\boldsymbol{B}||_F^2/C$$

Aside: Minimizing CE over $\boldsymbol{p}$ (or over a quantity on which $\boldsymbol{p}$ depends) is the same as minimizing the KL-divergence, $KL(\boldsymbol{q}, \boldsymbol{p}) = H(\boldsymbol{q}) + CE(\boldsymbol{q}, \boldsymbol{p})$ over it. The reason is $H(\boldsymbol{q})$ does not depend on $\boldsymbol{p}$.

## B.3 One-vs-rest (OVR) or One-vs-all use of two-class logistic regression**

There are two ways to deal with multiple classes using Log Reg. Say we have $K$ classes, $k = 0, 1, \ldots, K-1$.

1. Option 1: use what I have above

2. Option 2: use two class log reg in one-vs-rest or OVR mode. Idea of OVR:

   - First learn a $\theta_0$ for separating classes into 0 and everything else (combine classes 1,2,3, ... K-1 into the "not-zero" class)

   - Then learn a $\theta_1$ for separating classes into 1 and everything else (combine classes 0,2,3, ... K-1 into the "not-one" class)

   - Then learn a $\theta_2$ for separating classes into 2 and everything else (combine classes 0,1,3, ... K-1 into the "not-two" class)

   - And so on.

   - Output of above set of steps will be $K$ regression vectors $\theta_0, \theta_1, \ldots \theta_{K-1}$

   - Classification on a query $\boldsymbol{x}$:
     $$\hat{y} = \arg \max_{k=0,1,2,\ldots,K-1} \theta_k^\top \boldsymbol{x}$$

3. More details: see the video below

4. Also see https://en.wikipedia.org/wiki/Multiclass_classification

Video on One-vs-Rest (or One-vs-All) use of two-class Log Reg:
Lecture 6.7 of
https://www.youtube.com/watch?v=PPLop4L2eGk&list=PLLssT5z_DsK-h9vYZkQkYNWcItqhlRJLN

# C  Dealing with Categorical and Real-Valued Data Jointly:  Generative/Bayesiean models

## C.1  Modeling Categorical data

In some problems, the different features may be "categorical" instead of binary, which means feature $j$ takes one of $K_j$ possible values. For example $\boldsymbol{x}_j$ could be color of the front door of a house in case of the house price example and we assume $K_j = 5$ possible colors for example. There is no ordering to which color is preferred, thus the integer labels are arbitrary; they do not have a numerical meaning.

The model would then require us to learn $\psi_{j,y,k} := \Pr(\boldsymbol{x}_j = k|y)$ for $k = 1, 2, \ldots K_j$, $y = 0, 1$, $j = 1, 2, \ldots, n$.

Our data model then is as follows

$$p(\boldsymbol{x}|y) = \prod_{j=1}^{n} \prod_{k=1}^{K_j} \psi_{j,y,k}^{[\boldsymbol{x}_j == k]}$$

where $[x == k]$ takes the value 1 if $x = k$ and zero otherwise (MATLAB notation). Thus for MLE we need to maximize

$$\prod_{i} \left( \prod_{j=1}^{n} \prod_{k=1}^{K_j} (\psi_{j,y^{(i)},k})^{[\boldsymbol{x}_j^{(i)} == k]} \phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}} \right)$$

s.t. sum to one constraints on all the probabilities.

Then MLE estimates are given by

$$\hat{\psi}_{j,0,k} =$$

counts the number of times $y^{(i)} = 0$ and $\boldsymbol{x}_j^{(i)} = k$ in the training data and divides this by the number of times $y^{(i)} = 0$ (number of points from class zero in the training data). Do this for every value of $k = 1, 2, \ldots, K_j$ and for every feature $\boldsymbol{x}_j$, $j = 1, 2, \ldots, n$.

Do the same for class label 1.

## C.2  Dealing with categorical and real-valued data together in the Generative Learning (a.ka. Bayesian) framework **

Often some of the features can be categorical and some of them can be real-valued. Once we impose the Naive Bayes assumption (conditioned on the class label, different features are independent), this is easy to deal with. For simplicity suppose that the first $r$ features are real-valued and the rest $n - r$ are categorical with $k_j$ categories. Also assume the real-valued feature follow a Gaussian distribution with mean $\mu 0$ or $\mu_1$ (depending on the class label) and covariance matrix $\Sigma$. By Naive Bayes, $\Sigma$ is diagonal.

Then, joint likelihood becomes

$$\prod_{i=1}^{m} \left( \left( \prod_{j=1}^{r} \mathcal{N}(\boldsymbol{x}_j^{(i)}; (\mu_{y^{(i)}})_j, \sigma_j^2) \cdot \prod_{j=r+1}^{n} \prod_{k=1}^{K_j} \psi_{j,y^{(i)},k}^{[\boldsymbol{x}_j^{(i)} == k]} \right) \cdot \phi^{y^{(i)}} (1-\phi)^{1-y^{(i)}} \right)$$

s.t. sum to one constraints on $\phi$ and on all the $\psi_{j,y,k}$'s.

The MLE estimates are computed as follows.

$$\hat{\phi} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{1}(y^{(i)} = 1)$$

For $j = 1, 2, \ldots, r$,

$$(\hat{\mu}_0)_j = \frac{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 0)(\boldsymbol{x}^{(i)})_j}{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 0)}$$

$$(\hat{\mu}_1)_j = \frac{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 1)(\boldsymbol{x}^{(i)})_j}{\sum_{i=1}^m \mathbf{1}(y^{(i)} = 1)}$$

$$\hat{\sigma}_j^2 = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{x}^{(i)} - \mu_{y^{(i)}})_j^2, \ \ j = 1, 2, \ldots, n$$

For $j = r + 1, r + 2, \ldots, n$,

$$\hat{\psi}_{j,0,k} =$$

counts the number of times $y^{(i)} = 0$ and $\boldsymbol{x}_j^{(i)} = k$ in the training data and divides this by the number of times $y^{(i)} = 0$ (number of points from class zero in the training data). Do this for every value of $k = 1, 2, \ldots, K_j$ and for every feature $\boldsymbol{x}_j$, $j = 1, 2, \ldots, n$.

# D    Spam Filter details

### D.0.1    Spam Filter: entries of $z$ count the number of times a word occurs

The simplest spam filter explained above is not using word-counts (the number of times a word occurs), but just checking whether a word is present or not. A better model is to consider a different feature vector $\boldsymbol{z}$ with $\boldsymbol{z}_j$ being the number of times word $j$ occurs in the vocabulary (or dictionary).

$d$: number of words in any email – to use this model, we can use blank-space as one "word".

$n$: number of words in dictionary

$m$: number of training emails

Treating each word occurrence as independent (this is the naive Bayes assumption, it is actually not true since certain pairs of words are much more likely occur together, but simplifies our modeling), and assuming as before that $\psi_{j,y} = \Pr(\text{word } j \text{ is present in the email}|y)$, the feature vector $\boldsymbol{z}$ would be modeled by what is called a "multinomial distribution" as follows. This assumes that there are a total of $d$ words in the sample, i.e., that $\sum_j \boldsymbol{z}_j = d$ for all training emails. *To use this model, we can use blank-space as one "word".*

$$p(\boldsymbol{z}|y) = \binom{d}{\boldsymbol{z}_1, \boldsymbol{z}_2, \ldots, \boldsymbol{z}_n} \prod_{j=1}^n \psi_{j,y}^{\boldsymbol{z}_j}$$

This is called the multinomial distribution with parameters $\psi_{j,y}$.

**MLE** The joint likelihood of $m$ independent training samples can be expressed as follows:

$$J(\theta) = \prod_{i=1}^m p(y^{(i)}) p(\boldsymbol{z}^{(i)}|y^{(i)}) = \prod_{i=1}^m \phi^{y^{(i)}} (1 - \phi)^{(1 - y^{(i)})} \binom{d}{\boldsymbol{z}_1^{(i)}, \boldsymbol{z}_2^{(i)}, \ldots, \boldsymbol{z}_n^{(i)}} \prod_{j=1}^n \psi_{j,y^{(i)}}^{\boldsymbol{z}_j^{(i)}}$$

Learning parameters: we need to maximize the above over $\phi, \psi_{j,y}$ subject to the constraints that

$$0 \leq \psi_{j,y} \leq 1, \ 0 \leq \phi \leq 1$$

It can be shown that we get

$$\hat{\psi}_{j,0} = \frac{\sum_{i=1}^m \boldsymbol{z}_j^{(i)} \mathbb{1}(y^{(i)} = 0)}{\sum_{j=1}^n \sum_{i=1}^m \boldsymbol{z}_j^{(i)} \mathbb{1}(y^{(i)} = 0)}$$

In words this is the total number of times word $j$ occurs in **all** emails that are not spam (for which $y^{(i)} = 0$) divided by the total number of words in **all** emails that are not spam. We can similarly compute

$$\hat{\psi}_{j,1}$$

### D.0.2 Spam filter: model 3

A third possible model is as follows. Let $j$ referred to the $j$-th word in the email and let $\boldsymbol{x}_j$ be the location of the $j$-th word in the dictionary. So then this becomes "categorical data". We explain below how to deal with categorical data. In this case the length of the feature vector can be different for different emails. It will be $n^{(i)}$ if the email has $n^{(i)}$ distinct words.

### D.0.3 More ideas and related problems

A more advanced document model models co-occurrences of words.

Similar ideas are also used to classify blogs or webpages into various categories as well.

# E  Deep Learning / Deep Neural Networks: basic idea and training

## E.1  Introduction

From Wikipedia `https://en.wikipedia.org/wiki/Deep_learning`

- An older definition: Deep learning (DL) is a class of ML algorithms that uses multiple layers to progressively extract higher-level features from the raw input.

  - For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human, such as digits, letters, or faces.

- Modern definition: Deep learning (DL) is the subset of ML methods based on artificial neural networks with representation learning.

  - The goal is to "automate" human learning processes that map a source (e.g., an image of dog) to a learned object label (dog) by attempting to model what human brain neurons do.
  - The adjective "deep" refers to the use of multiple layers in the network.

- Most common DL methods involve supervised learning. Some newer approaches are also semi-supervised or unsupervised. Pros and Cons:

  - Supervised DL: training is very compute intensive and memory intensive. Query/test data processing is quick (linear in layer depth and width)
  - Unsupervised DL: no training; but then query/test data processing is extremely slow. Also needs a lot of memory (even more than supervised case).

- Types of architectures:

  - Feed-forward: multilayer perceptron (MLP) also called fully-connected net; Convolutional neural networks (CNN); more
  - Feed-back: Recurrent NN

- History taken from Wikipedia `https://en.wikipedia.org/wiki/Deep_learning`

  - The first deep learning multilayer perceptron (MLP) trained by stochastic gradient descent (Robbins, Munro 1951) was published in 1967 by Shun'ichi Amari.
  - In 1970, Seppo Linnainmaa published the basic idea of the algorithm that is now called "back-propoagation". Basic idea is to use chain rule of differentiation to compute the gradient of the cost function (energy) w.r.t. the weights of all layers.

* Details: In 1970, Seppo Linnainmaa published the reverse mode of automatic differentiation of discrete connected networks of nested differentiable functions. This became known as backpropagation. It is an efficient application of the chain rule derived by Gottfried Wilhelm Leibniz in 1673 to networks of differentiable nodes. More details: The terminology "back-propagating errors" was actually introduced in 1962 by Rosenblatt,[34][31] but he did not know how to implement this, although Henry J. Kelley had a continuous precursor of back-propagation[46] already in 1960 in the context of control theory.[31] In 1982, Paul Werbos applied backpropagation to MLPs in the way that has become standard.[47][48][31] In 1985, David E. Rumelhart et al. published an experimental analysis of the technique.[49]

- Convolutional neural networks (CNNs) with convolutional layers and downsampling layers began with the Neocognitron introduced by Kunihiko Fukushima in 1980.[50] In 1969, he also introduced the ReLU (rectified linear unit) activation function.

- The term Deep Learning was introduced to the ML community by Rina Dechter in 1986

## E.2 Visuals

`https://builtin.com/machine-learning/fully-connected-layer` – Fully connected (FC) vs Convolutional (CN) layer visual

`https://medium.com/@alejandro.itoaramendia/convolutional-neural-networks-cnns-a-complete-guide-` an example of a CNN - CN + Pool layers followed by FC

## E.3 Multi-layer Perceptron (MLP) or Fully Connected Feed-forward Network

The simplest neural net is the Feed-forward Network also called the Multilayer Perceptron or MLP. Each neuron receives a weighted sum of the outputs of the neurons of the previous layer, and applies a nonlinear "activation function" on this. Thus neuron $j$ in layer $k$ receives $o^{k-1}$ as input. It then computes

$$z_j^k = \sum_{i=1}^{r_k} w_{ij}^k o_i^{k-1}$$

and then outputs

$$o_j^k = g(z_j^k)$$

Here $g(z)$ is an element-wise nonlinearity. It could be the sigmoid function $\frac{1}{1+e^{-z}}$ or the Rectified Linear Unit (ReLU) function $\max(z, 0)$ or the tanh function.

Vectorizing the above, it can be expressed as

$$\boldsymbol{o}^k = g_{vec}(\boldsymbol{z}^k), \quad \boldsymbol{z}^k = \boldsymbol{W}^k \boldsymbol{o}^{k-1}$$

or equivalently,

$$\boldsymbol{z}^k = \boldsymbol{W}^k g_{vec}(\boldsymbol{z}^{k-1})$$

Here $g_{vec}(\boldsymbol{z})$ applies $g(.)$ to each entry of the vector $\boldsymbol{z}$.

The first layer takes the input $\boldsymbol{x}$ as input thus

$$\boldsymbol{z}^1 = \boldsymbol{x}$$

and computes $g_{vec}(\boldsymbol{z})$. Suppose the NN has 10 layers. The final (output) layer has only one neuron which outputs

$$\hat{y} = g(\boldsymbol{z}_1^{10})$$

Thus, the NN is

$$\hat{y} = g(\boldsymbol{W}^{10} g_{vec}(\boldsymbol{W}^9 g_{vec}(\boldsymbol{W}^8 \dots g_{vec}(\boldsymbol{x}))))$$

where $\boldsymbol{W}^{10}$ is a row vector (instead of a matrix).

Energy (cost function) for training is

$$E(\boldsymbol{W}^1, \boldsymbol{W}^2, \dots \boldsymbol{W}^{10}) := \sum_{i=1}^{m} (\boldsymbol{y}^{(i)} - \hat{y}^{(i)})^2 = \sum_{i=1}^{m} (\boldsymbol{y}^{(i)} - g(\boldsymbol{W}^{10} g_{vec}(\boldsymbol{W}^9 g_{vec}(\boldsymbol{W}^8 \dots g_{vec}(\boldsymbol{x}^{(i)})))))^2$$

## E.4 Training an MLP: Back-propagation **

Given training data $(\boldsymbol{x}^{(i)}, y^{(i)})$, $i = 1, 2, \dots, m$, we use gradient descent or stochastic / mini-batch GD to train. For all of these, the first task is to define the cost function $E(\hat{y}(\boldsymbol{x}), y)$ and to compute its gradient w.r.t. to each weight in each layer, i.e., compute

$$\frac{\partial E}{\partial w_{ij}^k}$$

This computation requires careful application of chain rule of differentiation. This leads to the following algorithm: consider a 10 layer NN and let $r_k$ denote the number of neurons in layer $k$

- For a given input $\boldsymbol{x}$, compute the outputs of all the layers and $\hat{y}$.

- Compute the following intermediate quantity:

$$\delta_i^k := \frac{\partial E}{\partial z_i^k}$$

using the following backward recursion (back-propagation)

  – compute

$$\delta_1^{10} = \frac{\partial E}{\partial \hat{y}}(\hat{y}, y) \cdot g'(z_1^{10})$$

  here $g'(z) = \frac{\partial g(z)}{\partial z}$

  – for each $k = 9, 8, \dots 1$, compute the following for each $i = 1, 2, \dots, r_k$:

$$\delta_i^k = g'(z_i^k) \sum_{j=1}^{r_{k+1}} w_{ij}^{k+1} \delta_j^{k+1}$$

  vectorized computation in MATLAB (or do similar in Python): $\boldsymbol{\delta}^k = g'(\boldsymbol{z}^k). * (\boldsymbol{W}^{k+1} \boldsymbol{\delta}^{k+1})$

- Compute

$$\frac{\partial E}{\partial w_{ij}^k} = \delta_j^k \cdot o_i^{k-1}$$

this can be vectorized too.

The above gives us $\frac{\partial E}{\partial w_{ij}^k}(\hat{y}(\boldsymbol{x}), y)$ for one input-output pair $\boldsymbol{x}, y$. The gradient w.r.t. the cost function that uses all the training data is thus

$$\frac{1}{m} \sum_{i=1}^{m} \frac{\partial E}{\partial w_{ij}^k}(\hat{y}(\boldsymbol{x}^{(i)}), y^{(i)})$$

## E.5 Understanding CNNs at a top level with visuals

## E.6 Aside: understanding convolution

See Dr. Zhengdao Wang's CNN notes (cnn-notes-zhengdao) `https://www.dropbox.com/s/4avs5jg3r5qctqe/cnn_notes_zhengdao.pdf?dl=0` for visuals and a quick review of LTI, 1D and 2D conv.

1D convolution:

2D convolution:

types of filters: low-pass, high-pass, band-pass,

low-pass filter, e.g., averaging filter: reduces noise but also blurs edges

high-pass filter, e.g., edge detector: nonzero response only at edges, zero response in parts of image that are constant (piecewise constant image)

high-pass in 2D : edges along different directions, can also be low-pass in one direction and high-pass in orthogonal direction

band-pass filters: created by applying high-pass filter on low-pass filter.

connection to wavelets.

## E.7 Convolutional Neural Network (CNN or ConvNet): basic idea

In case of MLP, each layer is fully-connected (FC), this means that all (or most) weights of the weight matrix are nonzero.

In case of CNN, this is not true. CNN consists of

- a convolutional (CN) layer followed by

- a pooling layer (pool),

- this sequence repeated a few times,

- finally 2-3 fully-connected (FC) layers.

Why each layer

- CN layer: extract useful features from images; most good features are local. Also much fewer parameters (due to paramater sharing – space-invariant filter for entire image)

- Pool layer: down-sample the image after the filtering to get a new smaller sized "image" (when you pass an image through any band-pass filter, e.g.m low pass or high pass filter, the amount of information reduces, consequently next layer should be smaller sized).

  How to downsample? Down-sample by 2: "low pass filter" the image followed by keeping every other pixel. For "low pass filter": either take average of neighbors or take max of neighbors. Usually max-pooling is popular. Averaging is linear space-invariant while Max is nonlinear and space-invariant.

- FC layer: CN and pooling layers help extract the features. FC layers take these features and provides a final classification output.

### E.7.1 Details of why CN + pooling instead of FC

- Local connectivity : useful when inputs are images

- Much fewer parameters (parameter sharing): $\ell_h^2 \cdot K'$, even with $\ell_h = 11, K' = 96$ (first layer of Alex Net), number of parameters is around 1000
  but if we vectorize even a $100 \times 100$ image and the next layer has more neurons than input layer (typically this is used), then there are more than 10000 parameters in an FC

- Heavily used in computer vision and image analyis applications.

## E.8 Different types of layer

Let $\boldsymbol{x}$ be the input to the layer and $\boldsymbol{z}$ the output.

### E.8.1 Fully connected (FC) layer:

$\boldsymbol{x}, \boldsymbol{z}$ are represented as vectors

- $\boldsymbol{x}$ is $n \times 1$, $\boldsymbol{z}$ is $n' \times 1$, $n'$ can be $>, =, < n$

- parameters: $\boldsymbol{W}$

- number of parameters: $n \cdot n'$

- apply nonlinearity, then weights' matrix

$$\boldsymbol{z}_{out} = \boldsymbol{W} g_{vec}(\boldsymbol{z}_{in})$$

### E.8.2 Convolutional (CN) layer:

$\boldsymbol{x}, \boldsymbol{z}$ are represented as 3D tensors

- $\boldsymbol{x}$ is $n_1 \times n_2 \times K$, $\boldsymbol{z}$ is $n'_1 \times n'_2 \times K'$ with $n'_1 \approx n_1$, $n'_2 \approx n_2$

- parameters: $h_k(u, v)$, $k = 1, 2, \ldots, K'$, $u = 0, 1, \ldots, \ell_h$, $v = 0, 1, \ldots, \ell_h$

- number of parameters: $\ell_h^2 \cdot K'$

- apply nonlinearity, convolution over the first 2 dimensions, sum over the third

$$\boldsymbol{z}(a, b, k) = \sum_w \left( \sum_u \sum_v h_k(u, v) g_{vec}(\ \boldsymbol{x}(a + u, b + v, w)\ ) \right)$$

- Details:

  - dealing with edges: zero pad, or reduce dimension of output
  - stride: above is stride $= 1$ (moving filter window by 1 pixel at a time) – default option
  - stride $= 2$, move filter window by 2 pixels
  - stride $= \ell_h$: non-overlapping window

### E.8.3 Pooling layer:

$\boldsymbol{x}, \boldsymbol{z}$ are represented as 3D tensors

- $\boldsymbol{x}, \boldsymbol{z}$ are represented as 3D tensors, $\boldsymbol{x}$ is $n_1 \times n_2 \times K$, $\boldsymbol{z}$ is $n'_1 \times n'_2 \times K$

- if using a $\ell \times \ell$ non-overlapping window for pooling, then $n'_1 = n_1/\ell, n'_2 = n_2/\ell$

- parameters: none

- number of parameters: zero

- applied on each $k = 1, 2, \ldots K$ "image" or "activation map" separately.

- helps reduce the output dimension (first 2 dimensions)

- max pooling over a $\ell \times \ell$ size window: keep the largest magnitude pixel value and throw the rest

- average pooling:

- stride $= \ell$: non-overlapping window – default option

## E.9   Other DL concepts

- Recurrent Neural Network (RNN)

- Generative Adversarial Networks GANs – used for DeepFake

- Auto-Encoder – use for denoising. Input and output is an image

- Unrolling

- Transformer

## E.10   Different NN architectures and when to use each

see the other handout NeuralNets-intro.pdf

## E.11   Reinforcement Learning

Use slides at `https://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture14.pdf`

- Example application: make a robot move.

  - Objective: Make the robot move forward
  - State: Angle and position of the joints
  - Action: Torques applied on joints
  - Reward: 1 at each time step that the robot is upright and moves forward

- There is an agent and the environment. We use $t$ to denote time.

- Agent or agent+environment has a state, $s_t$.

- Agent takes action $a_t$ based on current state $s_t$

  - Actions selected from a certain set. Example for robot it could be move front back left right one step.
  - Action selected using a policy $\pi$ that depends on $s_t$. This can be deterministic or random.

- The next state $s_{t+1}$, depends on $s_t$ and $a_t$

  - There is some randomness in this. This is specified by a probability transition matrix.

- Environment gives a reward $r_t$ to agent that depends on $s_t$ and $a_t$

  - There can be randomness in this.

- Agent receives

- ML Goal: design a "policy" that agent follows. Design this policy to maximize expected sum of future rewards (cumulative discounted reward)

- This framework is called Markov Decision Process or MDP.

  - Markov property: Current state completely characterises the state of the world

Details

- The value function for a given policy at state s, is the expected cumulative reward from following the policy from state s

- The Q-value function at state s and action a, is the expected cumulative reward from taking action a in state s and then following the policy

- The optimal Q-value function Q* is the maximum expected cumulative reward achievable from a given (state, action) pair.

- Goal is to find the optimal policy that helps achieve Q*.

- Two types of approaches:

  - Find or approximate $Q*$ first. Then use it to decide the best policy: given current state, find the action.

    * Q-learning: Use a function approximator to estimate the action-value function
    * If the function approximator is a deep neural network =¿ deep q-learning!
    * See slide 36 of above notes

  - Directly find a policy that optimizes the expected cumulative reward

    * Q function can be too complicated. But the policy can be much simpler
    * See slide 64 of above notes

# Final Exam - Spring 2025 EE 4250

- In-class. 2 sheets (4 sides) of HANDWRITTEN notes allowed.

  - 1 or 2 problems on Probability or Linear Algebra (midterm material)
  - 3 or 4 problems on ML algorithms
  - Last year exam discuss and in review next three lectures

- Topics on the Exam

  - Probability and Linear Algebra -
    * Probability Notes ML-cs229-prob.pdf (CourseNotes Module in Canvas)
    * Linear Algebra Notes Sections 2 and 3 of ML-cs229-linalg (CourseNotes Module in Canvas)
    * Linear Algebra Notes SVD (CourseNotes Module in Canvas)
    * Theory HWs 1-5
  - ML algorithms
    * ML-algorithms file (CourseNotes Module in Canvas)
  - ML Topics from ML algorithms file – 3-5 problems.
    * Monte Carlo methods and Data simulation
    * Gradient Descent (GD) basic GD only
    * Altnernating Minimization (AltMin)
    * Linear Regression
    * Logistic Regression - two class case only
    * GDA
    * PCA
    * Clustering and k-means clustering
    * Least Squares

- Topics that are good to know but will not be on the exam

  - Optimization overview
  - Multi-class Log Reg
  - Regularization - practical idea
  - Deep learning