

## Basics of Statistics

	Definition	Properties						
Sample Mean	A random variable. It can be an unbiased estimator of the real (population) mean. $\bar{X}_n = \frac{1}{n} \sum_i^n X_i$	$E[\bar{X}_n] = \mu$ $\text{var}(\bar{X}_n) = \frac{\sigma^2}{n} = (\text{SE})^2$ , where $\sigma^2$ is the variance of $X_i$						
Central Limit Theorem	Let $X_1, \dots, X_n$ form a random sample of size $n$ from a distribution with finite mean and variance, then for any fixed number $z$ : $\lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - \mu}{\sqrt{n}} \leq z\right) = \Phi(z)$	Interpreting: When $n$ is large, $\bar{X}_n$ is approximately normal: $\bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$ .						
Estimation	A guess of an unknown constant/parameter $\theta$ .  Estimator: a random variable. Estimate: a number, the realization of the estimator. This course use $\hat{\theta}$ to stand for both.	<b>Unbiased:</b> An estimator is unbiased if $E[\hat{\theta}] = \theta$ .  <b>Efficient:</b> An estimator is more efficient if it has a less variance.  <b>Consistent:</b> An estimator is consistent when it converges to $\theta$ in probability as $n \rightarrow \infty$ .  <b>Robust:</b> An estimator is robust if it still has a low bias when some assumptions of distribution are wrong.						
Sample Variance	$S^2 = \frac{1}{n-1} \sum_i^n (X_i - \bar{X}_n)^2$  It's a unbiased estimator of real (population) variance.							
Moment	Population Moments: $E[X], E[X^2], E[X^3], \dots$  Sample Moments: $\frac{1}{n} \sum_i^n X_i, \frac{1}{n} \sum_i^n X_i^2, \frac{1}{n} \sum_i^n X_i^3, \dots$							
? Method of moment								
Maximum Likelihood	This estimator of $\theta$ describes a PDFs which is most likely to have produced our data.  Let $L(\theta X) = \prod L(X_i; \theta)$ be the likelihood function. Then we can find $\hat{\theta}$ by maximize the $L(\theta X)$ (derivative=0).	If the maximization is difficult, try: 1. The log of the likelihood function. 2. Any monotone [increasing] transformation of the likelihood function.  Properties of maximum likelihood estimator: 1. If there is an efficient estimator, MLE will provide it. 2. Under certain regularity conditions, MLEs will have asymptotically normal distributions.						
Standard Error	$SE(\bar{X}_n) = \frac{\sigma}{\sqrt{n}} = \sqrt{\text{var}(\bar{X}_n)}$  $SE(\bar{X}_n) = \frac{\sigma}{\sqrt{n}}$ , when $\sigma$ is unknown, you need to estimate it.							
$\chi^2$ (Chi Square)	If $X_i$ is normal, then $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$ Where $S^2$ is the sample variance, $\sigma^2$ is the variance of $X_i$ , and $n-1$ is the degrees of freedom	Applications: - You think the true variance is $\sigma^2$ , and you want to use your data to make inferences about it.						
Student T	If $X \sim N(0, 1)$ and $Z \sim \chi^2_m$ and they're independent, then $\frac{X}{\sqrt{Z/m}} \sim t_n$ .	The t distribution is similar to the normal but has "fatter tails." It converges to the normal as $n \rightarrow \infty$ .						
F	If $X \sim \chi^2_n$ and $Z \sim \chi^2_m$ and they're independent, then $\frac{X}{m} \sim F_{n,m}$	Applications: - Compare the variances in two independent population. - check if 2 random samples are in the same population.						
Confidence Interval	A $1-\alpha$ confidence interval is an interval $[A(X_1, \dots, X_n), B(X_1, \dots, X_n)]$ such that: $P(A(X_1, \dots, X_n) \leq \theta \leq B(X_1, \dots, X_n)) = 1 - \alpha$	Case 1: sampling from normal, variance is known, find CI for the mean $P\left(\bar{X} + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(\alpha/2) \leq \mu \leq \bar{X} - \frac{\sigma}{\sqrt{n}} \Phi^{-1}(\alpha/2)\right) = 1 - \alpha$  Case 2: sampling from normal, variance is unknown, find CI for the mean $P\left(\bar{X} + \frac{S}{\sqrt{n}} t_{n-1}^{-1}(\alpha/2) \leq \mu \leq \bar{X} - \frac{S}{\sqrt{n}} t_{n-1}^{-1}(\alpha/2)\right) = 1 - \alpha$						
Hypothesis Testing	An hypothesis is an assumption about the distribution of a random variable in a population. A maintained hypothesis is one that cannot or will not be tested. A testable hypothesis is one that can be tested using evidence from a random sample.  The null hypothesis, $H_0$ , is the one that will be tested. The alternative hypothesis, $H_A$ , is a possibility (or series of possibilities) other than the null.  A simple hypothesis is one characterized by a single point, i.e., $\theta = 0$ . A composite hypothesis is one characterized by multiple points, i.e., $\theta > 0$ .  Critical region of the test, $C$ or $C_X$ , is the region of the support of the random sample for which we reject the null.  Test statistic $T$ is the random variable you calculated from the sample data, and you use it to calculate the p-value.	Testing steps: 1. You have a hypothesis (null hypothesis) 2. You observed something relevant (you got a sample) 3. You setup a threshold, i.e., 5% ( <b>significance level</b> ) 4. Calculate that ( <b>p-value</b> ) if the hypothesis is true, what is the probability (with the tail) that your sample will show up 5. If the p-value is less than your threshold (which means it's very unlikely to happen), then you should reject it						
Errors	<table border="1"><tr><td><math>H_0</math> true</td><td><math>H_0</math> false</td></tr><tr><td>accept <math>H_0</math></td><td>No error</td></tr><tr><td>reject <math>H_0</math></td><td>Type I error</td></tr></table> Significance Level = The Probability of Type I Error = $\alpha$ Confidence Level = $1 - \alpha$ Operating Characteristic = The Probability of Type II Error = $\beta$ Power = $1 - \beta$	$H_0$ true	$H_0$ false	accept $H_0$	No error	reject $H_0$	Type I error	$f_{\bar{X}}$ under $H_0$ vs $f_{\bar{X}}$ under $H_A$ 
$H_0$ true	$H_0$ false							
accept $H_0$	No error							
reject $H_0$	Type I error							

## Randomised Controlled Trials (RCTs)

Some Basic Concepts	<b>Causal effect for person i:</b> $Y_i(1) - Y_i(0)$ , 1 means being treated; only one of them can be observed.  <b>Treatment effect on the treated:</b> $E[Y_i(1) W_i = 1] - E[Y_i(0) W_i = 1]$ <b>Selection bias:</b> $E[Y_i(0) W_i = 1] - E[Y_i(0) W_i = 0]$	<b>Types of RCT:</b> <ul style="list-style-type: none"> <li>Completely randomized</li> <li>Stratified randomization</li> <li>Pairwise randomization</li> <li>Clustered randomization</li> </ul>
SUTVA	The potential outcomes for any unit do not vary with the treatments assigned to other units; and, for each unit, there are no different forms or versions of each treatment unit leading to different outcomes.	SUTVA wikipedia version: "the [potential outcome] observation on one unit should be unaffected by the particular assignment of treatments to the other units" (Cox 1958, §2.4)
ATE & Estimator	<b>Average treatment effect:</b> $E[Y^{obs} W_i = 1] - E[Y^{obs} W_i = 0]$  <b>Estimator:</b> $\hat{t} = \bar{Y}_t^{obs} - \bar{Y}_c^{obs}$  <b>Variance of estimator ((SE)<sup>2</sup>):</b> $V(\hat{t}) = \frac{S_t^2}{N_t} + \frac{S_c^2}{N_c}$  <b>Estimator of variance ((SE)<sup>2</sup>):</b> $\hat{V}(\hat{t}) = \frac{S_t^2}{N_t} + \frac{S_c^2}{N_c}$ $S_t^2 = \frac{1}{N_t - 1} \sum_{i \in W_t=1} (Y_i - \bar{Y}_t^{obs})^2$ $S_c^2 = \frac{1}{N_c - 1} \sum_{i \in W_c=0} (Y_i - \bar{Y}_c^{obs})^2$	
Confidence Interval of ATE Estimator	<b>Test statistic (t statistic):</b> $t = \frac{\hat{t}}{\sqrt{\hat{V}(\hat{t})}}$ , which follows a t-distribution.  <b>Confidence interval:</b> $CI_{1-\alpha} = (\hat{t} - t_{crit} * \sqrt{\hat{V}}, \hat{t} + t_{crit} * \sqrt{\hat{V}})$	<b>Small samples:</b> <ul style="list-style-type: none"> <li>find <math>t_{crit}</math> in t-table with <math>\alpha</math> and <math>df = N_t + N_c - 1</math></li> </ul>
Hypothesis Testing	$H_0$ : ate is 0 $H_1$ : ate is not 0  <b>Test statistic (t statistic):</b> $t = \frac{\bar{Y}_t^{obs} - \bar{Y}_c^{obs}}{\sqrt{\hat{V}(\hat{t})}}$	<b>Large samples:</b> <ul style="list-style-type: none"> <li>use normal approximation</li> <li><math>\alpha = 0.1</math></li> <li><math>\alpha = 0.05</math></li> </ul>
<b>Fisher Exact Test</b>		<b>Small samples:</b> <ul style="list-style-type: none"> <li><math>t</math> follows a t-distribution</li> <li><math>df = N - 1</math></li> <li><math>p = 2 * (1 - F_T(t))</math></li> </ul>
Hypothesis Testing	$H_0$ : $Y_i(1) = Y_i(0)$ for all $i$ $p = \Pr( T^{ave}(W, Y^{obs})  \geq  T^{ave}(W^{obs}, Y^{obs}) )$	<b>Steps:</b> <ol style="list-style-type: none"> <li>Calculate the estimate of ATE, which is <math>\hat{t} = \bar{Y}_t^{obs} - \bar{Y}_c^{obs}</math></li> <li>Get all possible permutations for treatments and controls</li> <li>Calculate <math>t^{ave} = \bar{Y}_t - \bar{Y}_c</math> in every permutation</li> <li><math>p = \frac{\text{count}(t^{ave} &gt; \hat{t})}{\text{count}(permutations)}</math></li> </ol>
Simulation	Instead of getting all permutations, we generate some random permutations, then p is: $p = \frac{\text{count}(t^{ave*} > \hat{t})}{\text{count}(random\ permutations)}$	
Some Important Concepts	<b>Exploratory data analysis:</b> In statistics, exploratory data analysis (EDA) is an approach analyzing data sets to summarize their main characteristics, often with visual methods. (from wikipedia)	
	<b>First order stochastic dominance (FOD) :</b> Distribution A to have first order stochastic dominance over Distribution B means: the cumulative distribution function (CDF) of Distribution A is always below or equal to the CDF of Distribution B, and must be strictly below at some value correct.	
K-S Test	<b>Kolmogorov-Smirnov test:</b> In statistics, the Kolmogorov-Smirnov test (K-S test or KS test) is a nonparametric test of the equality of continuous, one-dimensional probability distributions that can be used to compare a sample with a reference probability distribution (one-sample K-S test), or to compare two samples (two-sample K-S test). (from wikipedia)	<b>Applications:</b> <ul style="list-style-type: none"> <li>Compare distributions (r.v.'s)</li> <li>Test if two distributions are the same, e.g., <math>Y(0)</math> and <math>Y(1)</math></li> </ul>
Statistics	<b>Hypothesis testing:</b> $H_0: F = G$ (KS statistics tends to 0) $H_a: F > G$  <b>KS statistics:</b> $D_{nm} = \max[F_n(x) - G_m(x)]$  <b>KS statistics follows the distribution of:</b> $\frac{1}{n+m} D_{nm}^2$	<b>Reject <math>H_0</math> if <math>D_{nm} &gt; \frac{C(\alpha)(nm)}{n+m}</math>.</b> Use "ks.test" in R.
Nonparametric Regression	<b>Nonparametric regression:</b> Nonparametric regression is a category of regression analysis in which the predictor does not take a predetermined form but is constructed according to information derived from the data. (from wikipedia)	
Kernel Regression	<b>Kernel regression:</b> Kernel regression is a non-parametric technique in statistics to estimate the conditional expectation of a random variable. The objective is to find a non-linear relation between a pair of random variables X and Y.  $E[Y X] = g(X) \rightarrow y = g(x) + \epsilon$ $E[Y X = x] = \frac{1}{f(x)} \int y f(x, y) dy$	<b>Choice of bandwidth:</b> <ul style="list-style-type: none"> <li>A large bandwidth will lead to more bias (and less variance).</li> <li>A small bandwidth will lead to more variance (and less bias)</li> </ul>
	Kernel estimator replace $f(x, y)$ and $f(x)$ by their empirical estimates: $\hat{g}(x) = \frac{1}{f(x)} \int y f(x, y) dy = \frac{\sum_{i=1}^n y_i K(\frac{x-x_i}{h})}{\sum_{i=1}^n K(\frac{x-x_i}{h})}$ Where $h$ is the bandwidth (a positive number), the kernel estimate of the density of $x$ ; And $K()$ is your choice of kernel function, a density.	<b>Choice of kernel:</b> <ul style="list-style-type: none"> <li>Histogram: <math>K(u) = 1/2</math> if <math> u  \leq 1</math>, <math>K(u) = 0</math> otherwise.</li> <li>Epanechnikov <math>K(u) = \frac{3}{4}(1-u^2)</math> if <math> u  \leq 1</math>, <math>K(u) = 0</math> otherwise</li> <li>Quartic <math>K(u) = (\frac{3}{4}(1-u^2))^2</math> if <math> u  \leq 1</math>, <math>K(u) = 0</math> otherwise</li> </ul>
Cross Validation	<b>Cross Validation:</b> Cross-validation, sometimes called rotation estimation, or out-of-sample testing is any of various similar model validation techniques for assessing how the results of a statistical analysis will generalize to an independent data set. It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice.	Optimize your choice: using cross-validation.

<p><b>Linear Model &amp; Linear Regression</b></p> <p><b>Linear model</b>  <math>y_i = \beta_0 + \beta_1 x_i + \epsilon_i, i = 1, \dots, n</math></p> <ul style="list-style-type: none"> <li>- <math>Y</math>: dependent variable</li> <li>- <math>X</math>: regressor/explanatory variable</li> <li>- <math>\beta_0, \beta_1</math>: regression coefficients</li> <li>- <math>\epsilon</math>: error</li> </ul> <p><b>Linear regression</b>  The way we estimate the parameters (<math>\beta_0</math> and <math>\beta_1</math>).</p> <p><b>Assumptions</b></p> <ol style="list-style-type: none"> <li><math>X_i, \epsilon_i</math> uncorrelated</li> <li>Identification: <math>\frac{1}{n} \sum_i (X_i - \bar{X})^2 &gt; 0</math></li> <li>Zero mean: <math>E[\epsilon_i] = 0</math></li> <li>Homoskedasticity: <math>\text{var}(\epsilon_i) = \sigma^2</math> for all <math>i</math></li> <li>No serial correlation: <math>E[\epsilon_i \epsilon_j] = 0, \text{ if } i \neq j</math></li> </ol>	<p><b>Linear Regression</b></p> <p>For:</p> <ul style="list-style-type: none"> <li>- Prediction</li> <li>- Determining causality</li> <li>- understanding the world better</li> </ul>	<p><b>Dummy Variables</b></p> <p><b>Dummy Variable</b>  Dummy variable (indicator variable) is one that takes the value 0 or 1 to indicate the absence or presence of some categorical effect that may be expected to shift the outcome. (from wikipedia)</p> <p><b>Categorical Variables:</b></p> <ul style="list-style-type: none"> <li>• Use multiple dummy variables to build categories</li> <li>• “-1” to avoid collinear problem</li> </ul>	<p><b>Example:</b>  <math>Y_i = \alpha + \beta D_i + \gamma M_i + \delta M_i * D_i + \epsilon_i</math></p> <p><math>D_i</math> takes the value 1 if the observation is in group A, and 0 if in group B.  <math>\hat{\beta} = \bar{Y}_A - \bar{Y}_B</math></p> <p>You can always estimate the difference between the treatment and control group for an RCT using an OLS regression framework. The standard errors will be slightly different from the Neyman standard errors, but it won't matter that much if the samples are large enough, and similar in treatment and control groups.</p>																																																						
<p><b>Least Squares Estimator</b></p> $\min_{\beta} \sum_i (Y_i - \beta_0 - \beta_1 X_i)^2$ $\hat{\beta}_1 = \frac{1}{n} \sum_i (X_i - \bar{X})(Y_i - \bar{Y})$ $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$ <p>Let <math>\bar{X} = \frac{1}{n} \sum X_i</math> and <math>\hat{\sigma}_x^2 = \frac{1}{n} \sum (X_i - \bar{X})^2</math>.</p> <table border="0"> <tr> <td style="text-align: center;"><math>\hat{\beta}_0</math></td> <td style="text-align: center;"><math>\hat{\beta}_1</math></td> <td style="text-align: center;"><math>\text{mean}</math></td> <td style="text-align: center;"><math>\text{variance}</math></td> <td style="text-align: center;"><math>\text{covariance}</math></td> </tr> <tr> <td><math>\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}</math></td> <td><math>\hat{\beta}_1 = \frac{1}{n} \sum (X_i - \bar{X}) \frac{\hat{\sigma}_x^2}{n \hat{\sigma}_x^2}</math></td> <td><math>\hat{\sigma}_x^2 = \frac{1}{n} \sum (X_i - \bar{X})^2</math></td> <td><math>\hat{\sigma}_y^2 = \frac{1}{n} \sum (Y_i - \bar{Y})^2</math></td> <td><math>\text{cov}(\beta_0, \beta_1) = -\hat{\sigma}_x^2 / n \hat{\sigma}_y^2 &lt; 0</math></td> </tr> </table> <p>Estimation of error variance <math>\sigma^2</math>:</p> $\hat{\sigma}^2 = \frac{1}{n-2} \sum_i \epsilon_i^2$	$\hat{\beta}_0$	$\hat{\beta}_1$	$\text{mean}$	$\text{variance}$	$\text{covariance}$	$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$	$\hat{\beta}_1 = \frac{1}{n} \sum (X_i - \bar{X}) \frac{\hat{\sigma}_x^2}{n \hat{\sigma}_x^2}$	$\hat{\sigma}_x^2 = \frac{1}{n} \sum (X_i - \bar{X})^2$	$\hat{\sigma}_y^2 = \frac{1}{n} \sum (Y_i - \bar{Y})^2$	$\text{cov}(\beta_0, \beta_1) = -\hat{\sigma}_x^2 / n \hat{\sigma}_y^2 < 0$	<p>Properties of estimator (in model):</p> <ul style="list-style-type: none"> <li>- Least variance (most efficient)</li> <li>- Unbiased</li> <li>- Normal errors mean normal estimates</li> </ul> <p><b>Some comparative statics:</b></p> <ul style="list-style-type: none"> <li>--A larger <math>\sigma^2</math> means larger <math>\text{Var}(\hat{\beta})</math></li> <li>--A larger <math>\hat{\sigma}_x^2</math> means smaller <math>\text{Var}(\hat{\beta})</math></li> <li>--A larger <math>n</math> means smaller <math>\text{Var}(\hat{\beta})</math></li> <li>--If <math>\bar{X} &gt; 0</math>, <math>\text{Cov}(\beta_0, \beta_1) &lt; 0</math></li> </ul> <p>If we use the stronger assumption that the errors are i.i.d. <math>N(0, \sigma^2)</math>, we obtain the result that <math>\beta_0</math> and <math>\beta_1</math> will also have normal distributions.</p>	<p><b>Difference in differences</b></p> $Y_i = \alpha + \beta D_i + \gamma M_i + \delta M_i * D_i + \epsilon_i$ <p><math>\alpha</math>: mean for women in the control group  <math>\beta</math>: treatment main effect  <math>\gamma</math>: gender main effect  <math>\delta</math>: difference between the treatment effect for males and for female, interaction effect</p> <p><b>DID table:</b></p> <table border="1"> <thead> <tr> <th></th> <th><math>D = 0</math></th> <th><math>D = 1</math></th> </tr> </thead> <tbody> <tr> <td><math>M = 0</math></td> <td><math>\alpha</math></td> <td><math>\alpha + \beta</math></td> </tr> <tr> <td><math>M = 1</math></td> <td><math>\alpha + \gamma</math></td> <td><math>\alpha + \beta + \gamma + \delta</math></td> </tr> </tbody> </table> $\delta = ((\alpha + \beta + \gamma + \delta) - (\alpha + \gamma)) - ((\alpha + \beta) - \alpha)$		$D = 0$	$D = 1$	$M = 0$	$\alpha$	$\alpha + \beta$	$M = 1$	$\alpha + \gamma$	$\alpha + \beta + \gamma + \delta$	<p>A variation example:</p> $S_{ijk} = c_1 + \alpha_{1j} + \beta_{1k} + (P_j * T_i) \gamma_{1j} + \epsilon_{ijk}, \quad (1)$ <p>where</p> <ul style="list-style-type: none"> <li>• <math>S_{ijk}</math> is the education of individual <math>i</math> born in region <math>j</math> in year <math>k</math>,</li> <li>• <math>T_i</math> is a dummy indicating whether the individual belongs to the “young” cohort in the subsample,</li> <li>• <math>P_j</math> denotes the intensity of the program in the region of birth (number of school built)</li> <li>• <math>c_1</math> is a constant,</li> <li>• <math>\beta_{1k}</math> is a set of cohort-of-birth fixed effects [in practice, a series of dummies=1 for each year of birth, omit 1]</li> <li>• <math>\alpha_{1j}</math> is a set of district-of-birth fixed effects [in practice, a series of dummies=1 for each district of birth, omit 1]</li> </ul>																																			
$\hat{\beta}_0$	$\hat{\beta}_1$	$\text{mean}$	$\text{variance}$	$\text{covariance}$																																																					
$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$	$\hat{\beta}_1 = \frac{1}{n} \sum (X_i - \bar{X}) \frac{\hat{\sigma}_x^2}{n \hat{\sigma}_x^2}$	$\hat{\sigma}_x^2 = \frac{1}{n} \sum (X_i - \bar{X})^2$	$\hat{\sigma}_y^2 = \frac{1}{n} \sum (Y_i - \bar{Y})^2$	$\text{cov}(\beta_0, \beta_1) = -\hat{\sigma}_x^2 / n \hat{\sigma}_y^2 < 0$																																																					
	$D = 0$	$D = 1$																																																							
$M = 0$	$\alpha$	$\alpha + \beta$																																																							
$M = 1$	$\alpha + \gamma$	$\alpha + \beta + \gamma + \delta$																																																							
<p><b>Multivariate Linear Model</b></p> <p><b>Multivariate Linear Model</b></p> <p><math>Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \epsilon_i, i = 1, \dots, n</math></p> <p>Or</p> <p><math>X\beta = Y</math></p> <p><math>Y = (Y_1, \dots, Y_n)^T</math>, a <math>n \times 1</math> column vector</p> <p><math>\epsilon = (\epsilon_1, \dots, \epsilon_n)^T</math>, a <math>n \times 1</math> column vector</p> <p><math>X = \begin{bmatrix} X_{01} &amp; \dots &amp; X_{0k} \\ X_{11} &amp; \dots &amp; X_{1k} \\ \vdots &amp; \ddots &amp; \vdots \\ X_{n1} &amp; \dots &amp; X_{nk} \end{bmatrix}</math>, a <math>n \times (k+1)</math> matrix (<math>X_{0i} = 1</math>, for <math>\beta_0</math>)</p> <p>Assumptions:</p> <ol style="list-style-type: none"> <li>1) <math>n &gt; k + 1</math></li> <li>2) <math>\epsilon \sim N(0, \sigma^2 I_n)</math>, where <math>I_n</math> is a <math>n \times n</math> identity matrix</li> </ol>	<p>Interpreting assumptions:</p> <ul style="list-style-type: none"> <li>- need to have more observations than regressors</li> <li>- can't have any regressors that do not have positive sample variation</li> <li>- can't have any regressors that are linear functions of one or more other regressors</li> </ul> <p><math>E[\beta] = \beta</math></p> <p><math>\text{cov}(\beta) = \sigma^2 (X^T X)^{-1}</math></p> <p><math>\hat{\sigma}^2 = \frac{\epsilon^T \epsilon}{n-k}</math></p>	<p><b>Transformations</b></p> $Y_i = X_{11}^{\beta_1} X_{21}^{\beta_2} \epsilon_i \Rightarrow \log(Y_i) = \beta_0 + \beta_1 \log X_{11} + \beta_2 \log X_{21} + \epsilon$ <p><b>Education formulation:</b></p> $\log Y_i = \beta_0 + \beta_1 S_i + \epsilon_i$ <p><b>Box Cox Transformation:</b></p> $Y_i = \frac{1}{\beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \epsilon_i} \Rightarrow \frac{1}{Y_i} = \beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \epsilon_i$ <p><b>Discrete choice model:</b></p> $P_i = \frac{e^{\beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \epsilon_i}}{1 + e^{\beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \epsilon_i}} \Rightarrow Y_i = \log \frac{P_i}{1 - P_i} = \beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \epsilon_i$	<p><b>Non-Linear Models</b></p> <p><b>Polynomial Models</b></p> $Y_i = \beta_0 + \beta_1 X_{11} + \beta_2 X_{11}^2 + \dots + \beta_k X_{11}^k + \epsilon_i$ <p><b>Keys:</b></p> <ul style="list-style-type: none"> <li>- You can chose series expansion, or orthogonal polynomials</li> <li>- If you assume the model is known, this is OLS</li> <li>- If you assume the model is unknown, this is a non-parametric method: series regression</li> </ul> <p><b>Using Dummies for Approximation</b></p> <ul style="list-style-type: none"> <li>• Partition the range of <math>X</math> is interval, <math>X_0, \dots, X_J</math></li> <li>• Define the dummies as:</li> <ul style="list-style-type: none"> <li><math>D_{1i} = I_{[X_0 \leq X_i &lt; X_1]}</math></li> <li><math>D_{2i} = I_{[X_1 \leq X_i &lt; X_2]}</math></li> <li><math>\vdots</math></li> <li><math>D_{Ji} = I_{[X_{J-1} \leq X_i &lt; X_J]}</math></li> </ul> <li>• you can run regression:</li> </ul> $Y_i = \beta_1 D_{1i} + \beta_2 D_{2i} + \dots + \beta_J D_{Ji} + \epsilon_i$																																																						
<p><b>Estimator</b></p> <p><math>\hat{\beta} = (X^T X)^{-1} X^T Y</math></p> <p><math>T = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})}</math> follows a t-distribution.</p> <p><b>Hypothesis testing</b></p> <p><math>H_0: \beta_i = c</math></p> <p><math>T = \frac{\hat{\beta}_i - c}{SE(\hat{\beta}_i)}</math>, where <math>SE(\hat{\beta}_i) = (\sigma^2 (X^T X)^{-1})_{ii}^{1/2}</math></p> <p><math>R</math> is a <math>1 \times (k+1)</math> matrix:</p> <p><math>H_0: R\beta = c</math></p> <p><math>T = \frac{R\hat{\beta} - c}{SE(R\hat{\beta})}</math>, where <math>SE(R\hat{\beta}) = (\sigma^2 R(X^T X)^{-1} R^T)^{1/2}</math></p>	<p>For the hypothesis <math>H_0: \beta_i = c</math> versus <math>H_0: \beta_i \neq c</math>, the F-test is equivalent to the t-test (<math>\sqrt{F}= T </math>).</p>	<p>Locally Linear Regression</p> <p>Like a kernel, but using a linear regression in each little interval instead.</p> <p><b>Regression Discontinuity Design</b></p> <p>A regression discontinuity design (RDD) is a quasi-experimental pretest-postest design that elicits the causal effects of interventions by assigning a cutoff or threshold above or below which an intervention is assigned. (from wikipedia)</p> <p><b>Running Variable</b></p> <p>RD is appropriate in any circumstance where some treatment shift discontinuously with a variable <math>a</math>, called the running variable.</p>	<p>Locally linear estimation would have no bias if the true model were linear.</p> <p>Locally linear estimation removes a bias term from the kernel estimator, that makes it have better behavior near the boundary of the x's and smaller MSE everywhere.</p> <p>"You cannot do a regression discontinuity analysis without a graph."</p> <p><b>An Example</b></p> <p>Simplest model: <math>Y_i = \beta_0 + \beta_1 D_{ai} + \beta_2 a_i + \epsilon_i</math> (same slopes for two parts)</p> <p>Improvement:</p> <ul style="list-style-type: none"> <li>- Add piecewise polynomials</li> <li>- Fit a polynomial on each side:</li> <ul style="list-style-type: none"> <li><math>Y_i = \beta_0 + \beta_1 D_{ai} + \beta_2 (a_i - a_0) + \beta_3 (a_i - a_0) * D_{ai} + \dots + \epsilon_i</math></li> </ul> </ul> <p><b>Omitted Variable Bias</b></p> <p>Omitted variable bias depends on:</p> <ol style="list-style-type: none"> <li>1. How important is <math>X_j</math> in the original model</li> <li>2. How correlated it is with <math>X_i</math></li> </ol> <p><b>Proof:</b></p> $\hat{\alpha}_1 = \frac{\text{Cov}(Y_i, X_{1i})}{V(X_{1i})}$ <p>substituting for <math>Y_i</math> we get:</p> $\frac{\text{Cov}(\beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \epsilon_i, X_{1i})}{V(X_{1i})}$ $= \frac{\beta_1 V(X_{1i}) + \beta_2 \text{Cov}(X_{2i}, X_{1i}) + \text{Cov}(\epsilon_i, X_{1i})}{V(X_{1i})}$ $= \beta_1 + \beta_2 \beta_1$ <p>Advanced techniques a reasonable causal model:</p> <ul style="list-style-type: none"> <li>- Matching methods</li> <li>- Propensity score matching</li> <li>- Machine Learning techniques</li> </ul>																																																						
<p><b>Testing Restrictions</b></p> <p><math>R</math> is a <math>r \times (k+1)</math> matrix of restrictions.</p> <p><math>H_0: R\beta = c</math></p> <p>Testing steps:</p> <ol style="list-style-type: none"> <li>1. Estimate the unrestricted model.</li> <li>2. Impose the restrictions of the null and estimate that model.</li> <li>3. Compare the goodness-of-fit of the models. If the restrictions don't really affect the fit of the model much, then the null is probably true or close to true.</li> </ol> <p><math>SSR_R - SSR_H</math></p> <p><math>T = \frac{r}{\frac{SSR_H - SSR_R}{SSR_H}}</math></p> <p><math>r</math> is the restriction number (the number of independent variables are dropped)</p> <p><math>n</math> is the observations number, <math>k+1</math> is the variables number</p> <p><math>T \sim F_{r, n-(k+1)}</math></p>	<p><b>R functions - part 1</b></p> <table border="1"> <thead> <tr> <th>name</th> <th>description</th> </tr> </thead> <tbody> <tr> <td><code>dname()</code></td> <td>density or probability function</td> </tr> <tr> <td><code>pname()</code></td> <td>cumulative density function</td> </tr> <tr> <td><code>qname()</code></td> <td>quantile function</td> </tr> <tr> <td><code>Rname()</code></td> <td>random deviates</td> </tr> </tbody> </table> <p><b>Restriction examples:</b></p> <p>If, for instance, <math>R = [0 \ 1 \ 0 \dots \ 0]</math> and <math>c = [0]</math> that corresponds to <math>H_0: \beta_1 = 0</math>.</p> <p>If, for instance, <math>R = [0 \ 1 \ 0 \dots \ 0]</math> and <math>c = [0]</math></p> $\begin{bmatrix} 0 \\ 0 \\ 0 \\ \dots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ \dots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ \dots \\ 0 \end{bmatrix}$ <p>that corresponds to <math>H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0</math>.</p> <p>If, for instance, <math>R = [0 \ 1 \ -1 \ \dots \ 0]</math> and <math>c = [0]</math></p> $\begin{bmatrix} 0 \\ 0 \\ 0 \\ \dots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ \dots \\ 0 \end{bmatrix} \begin{bmatrix} 5 \\ \dots \\ -2 \end{bmatrix}$ <p>that corresponds to <math>H_0: \beta_1 = \beta_2, \beta_3 = 5</math>, and <math>\beta_k = -2</math>.</p>	name	description	<code>dname()</code>	density or probability function	<code>pname()</code>	cumulative density function	<code>qname()</code>	quantile function	<code>Rname()</code>	random deviates	<p><b>Omitted Variable Bias</b></p> <p><b>OVB</b></p> <p>In statistics, omitted-variable bias (OVB) occurs when a statistical model leaves out one or more relevant variables. The bias results in the model attributing the effect of the missing variables to the estimated effects of the included variables. (from wikipedia)</p> <p>Correct model: <math>Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \epsilon_i</math></p> <p>Estimated model: <math>Y_i = \alpha_0 + \alpha_1 X_{1i} + w_i</math></p> <p>Ancillary regression: <math>X_{2i} = \delta_0 + \delta_1 X_{1i} + \xi_i</math></p> <p>OVB: <math>OVB = \alpha_1 - \beta_1 = \delta_1 \beta_1</math></p> <p><b>R functions - part 2</b></p> <table border="1"> <thead> <tr> <th>distribution</th> <th>R name</th> <th>distribution</th> <th>R name</th> </tr> </thead> <tbody> <tr> <td>Beta</td> <td>beta</td> <td>Lognormal</td> <td>lnorm</td> </tr> <tr> <td>Binomial</td> <td>binom</td> <td>Negative Binomial</td> <td>nbinom</td> </tr> <tr> <td>Cauchy</td> <td>cauchy</td> <td>Normal</td> <td>norm</td> </tr> <tr> <td>Chisquare</td> <td>chisq</td> <td>Poisson</td> <td>pois</td> </tr> <tr> <td>Exponential</td> <td>exp</td> <td>Student t</td> <td>t</td> </tr> <tr> <td>F</td> <td>f</td> <td>Uniform</td> <td>unif</td> </tr> <tr> <td>Gamma</td> <td>gamma</td> <td>Tukey</td> <td>tukey</td> </tr> <tr> <td>Geometric</td> <td>geom</td> <td>Weibull</td> <td>weib</td> </tr> <tr> <td>Hypergeometric</td> <td>hyper</td> <td>Wilcoxon</td> <td>wilcox</td> </tr> <tr> <td>Logistic</td> <td>logis</td> <td></td> <td></td> </tr> </tbody> </table>	distribution	R name	distribution	R name	Beta	beta	Lognormal	lnorm	Binomial	binom	Negative Binomial	nbinom	Cauchy	cauchy	Normal	norm	Chisquare	chisq	Poisson	pois	Exponential	exp	Student t	t	F	f	Uniform	unif	Gamma	gamma	Tukey	tukey	Geometric	geom	Weibull	weib	Hypergeometric	hyper	Wilcoxon	wilcox	Logistic	logis			
name	description																																																								
<code>dname()</code>	density or probability function																																																								
<code>pname()</code>	cumulative density function																																																								
<code>qname()</code>	quantile function																																																								
<code>Rname()</code>	random deviates																																																								
distribution	R name	distribution	R name																																																						
Beta	beta	Lognormal	lnorm																																																						
Binomial	binom	Negative Binomial	nbinom																																																						
Cauchy	cauchy	Normal	norm																																																						
Chisquare	chisq	Poisson	pois																																																						
Exponential	exp	Student t	t																																																						
F	f	Uniform	unif																																																						
Gamma	gamma	Tukey	tukey																																																						
Geometric	geom	Weibull	weib																																																						
Hypergeometric	hyper	Wilcoxon	wilcox																																																						
Logistic	logis																																																								

## 6.86x Machine Learning with Python

This is a cheat sheet for machine learning based on the online course given by Prof. Tommi Jaakkola and Prof. Regina Barzilay. Compiled by Janus B. Advinencia.

Last Updated January 17, 2020

### Linear Classifiers

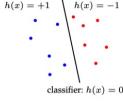
#### Introduction to Machine Learning

What is machine learning? Machine learning as a discipline aims to design, understand and apply computer programs that learn from experience (i.e., data) for the purpose of modeling, prediction, or control.

#### Types of Machine Learning

- Supervised learning: prediction based on examples of correct behavior
- Unsupervised learning: no explicit target, only data, goal is to model/discover
- Semi-supervised learning: supplement limited annotations with unsupervised learning
- Active learning: learn to query the examples actually needed for learning
- Transfer learning: how to apply what you have learned from A to B
- Reinforcement learning: learning to act, not just predict; goal is to optimize the consequences of actions

#### Linear Classifier and Perceptron



#### Key Concepts

- feature vectors, labels:  $x \in \mathbb{R}^d, y \in \{-1, +1\}$
- training set:  $S_n = \{(x^{(i)}, y^{(i)}), i=1, \dots, n\}$
- classifier:  $h: \mathbb{R}^d \rightarrow \{-1, +1\}$   
 $X^n = \{x \in \mathbb{R}^d : h(x) = +1\}$   
 $X^- = \{x \in \mathbb{R}^d : h(x) = -1\}$
- training error:  $E_n(h) = \frac{1}{n} \sum_i \mathbb{I}[(x^{(i)}) \neq y^{(i)}]$   
 $\llbracket h(x^{(i)}) \neq y^{(i)} \rrbracket = \begin{cases} 1 & \text{if error} \\ 0 & \text{otherwise} \end{cases}$
- test error:  $\mathcal{E}(h)$

**Collaborative Filtering** Our goal is to come up with a matrix  $X$  that has no blank entries and whose  $(a, i)$ -th entry  $X_{ai}$  is a user rating of movie  $i$ . Let  $D$  be the set of all  $(a, i)$ 's for which a user rating  $Y_{ai}$  exists. A naive approach is to minimize the objective function

$$J(X) = \sum_{(a,i) \in D} \frac{1}{2} (Y_{ai} - X_{ai})^2 + \frac{\lambda}{2} \sum_{i,k} X_{ik}^2.$$

The results are

$$\hat{X}_{ai} = \frac{Y_{ai}}{1 + \lambda} \quad \text{for } (a, i) \in D$$

$$\hat{X}_{ai} = 0 \quad \text{for } (a, i) \notin D.$$

The problem with this approach is that there is no connection between the entries of  $X$ . We can impose additional constraint on  $X$ :

$$X = UV^T$$

for some  $n \times k$  matrix  $U$  and  $d \times m$  matrix  $V^T$ , where  $d$  is the rank of the matrix  $X$ .

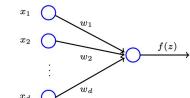
**Alternating Minimization** Assume that  $U$  and  $V$  are rank  $k$  matrices. Then, we can write the objective function as

$$J(X) = \sum_{(a,i) \in D} \frac{1}{2} (Y_{ai} - u_a v_i)^2 + \frac{\lambda}{2} \sum_k u_k^2.$$

### Neural Networks

#### Introduction to Feedforward Neural Networks

A Unit in a Neural Network A neural network unit is a primitive neural network that consists of only the input layer, and an output layer with only one output.



A neural network unit computes a non-linear weighted combination of its input:

$$g = f(z) \quad \text{where } z = w_0 + \sum_{i=1}^d x_i w_i$$

where  $w_i$  are the **weights**,  $z$  is a number and is the weighted sum of the inputs  $x_i$ , and  $f$  is generally a non-linear function called the **activation function**.

**Linear Function**  $f(z) = z$

**Rectified Linear Unit (ReLU)**  $f(z) = \max\{0, z\}$

• set of classifiers:  $h \in \mathcal{H}$

**Linear Classifiers through the Origin** We consider functions of the form  $h(x; \theta) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d) = \text{sign}(\theta \cdot x)$ .



**Linear Classifiers with Offset** We can consider functions of the form  $h(x; \theta) = \text{sign}(\theta \cdot x + \theta_0)$ .

where  $\theta_0$  is the offset parameter.

**Linear Separation** Training examples  $S_n$  are **linearly separable** if there exists a parameter vector  $\theta$  and offset parameter  $\theta_0$  such that  $y^{(i)} (\theta \cdot x^{(i)} + \theta_0) > 0$  for all  $i = 1, \dots, n$ .

**Training Error** The training error for a linear classifier is

$$E_n(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y^{(i)} (\theta \cdot x^{(i)} + \theta_0) \leq 0].$$

**Perceptron Algorithm**

```
procedure PERCEPTRON({(x^{(i)}, y^{(i)}), i=1, ..., n}, T)
    θ = 0 (vector)
    for t = 1, ..., T do
        for i = 1, ..., n do
            if y^{(i)} (θ · x^{(i)}) ≤ 0 then
                θ = θ + y^{(i)} x^{(i)}
        return θ
```

**Perceptron Algorithm (with offset)**

```
procedure PERCEPTRON({(x^{(i)}, y^{(i)}), i=1, ..., n}, T)
    θ = 0 (vector)
    for t = 1, ..., T do
        for i = 1, ..., n do
            if y^{(i)} (θ · x^{(i)} + θ_0) ≤ 0 then
                θ = θ + y^{(i)} x^{(i)}
                θ_0 = θ_0 + y^{(i)}
        return θ, θ_0
```

**Convergence Assumption**

- There exists  $\theta^*$  such that  $\frac{y^{(i)} (\theta^* \cdot x^{(i)})}{\|x^{(i)}\|} \geq \gamma$  for all  $i = 1, \dots, n$  for some  $\gamma > 0$ .
- All examples are bounded  $\|x^{(i)}\| \leq R, i = 1, \dots, n$ .

Then the number  $k$  of updates made by the perceptron algorithm is bounded by  $R^2 / \gamma^2$ .

#### Hinge Loss, Margin Boundaries and Regularization

**Distance from a Line to a Point** The perpendicular distance from a line with equation  $\theta \cdot x + \theta_0 = 0$  to a point with coordinates  $x_0$  is

$$d = \frac{|\theta \cdot x_0 + \theta_0|}{\|\theta\|}$$

$\eta_t$  is the learning rate which can vary at every iteration.

**Support Vector Machine**

- Support Vector Machine finds the maximum margin linear separator by solving the quadratic program that corresponds to  $J(\theta, \theta_0)$ .
- In the realizable case, if we disallow any margin violations, the quadratic program we have to solve is

Find  $\theta, \theta_0$  to minimize  $\frac{1}{2} \|\theta\|^2$  subject to

$$y^{(i)} (\theta \cdot x^{(i)} + \theta_0) \geq 1, \quad i = 1, \dots, n$$

#### Nonlinear Classification, Linear Regression, Collaborative Filtering

##### Linear Regression

**Empirical Risk**

$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \theta \cdot x^{(i)})^2 \quad \text{squared error}$$

**Gradient-based Approach** We can use stochastic gradient descent to find the minima of the empirical risk.

```
Algorithm Initialize θ = 0.
Randomly pick i = {1, ..., n}.
θ = θ + η (y^{(i)} - θ · x^{(i)}) x^{(i)}.
```

$\eta_t$  is the learning rate.

**Closed Form Solution**

$$\text{Let } A = \frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^T \quad \text{and } B = \frac{1}{n} \sum_{i=1}^n y^{(i)} x^{(i)}.$$

Then,

$$\hat{\theta} = A^{-1} B.$$

In matrix notation, this is

$$\hat{\theta} = (X^T X)^{-1} X^T Y.$$

**Generalization and Regularization**

**Ridge Regression**: The loss function is

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n [\text{Loss}_i(y^{(i)}, \theta \cdot x^{(i)} + \theta_0)] + \frac{\lambda}{2} \|\theta\|^2$$

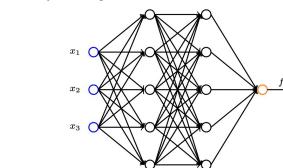
$\lambda$  is the regularization factor.

where  $\lambda$  is the regularization factor. We can find its minima using gradient-based approach.

```
Algorithm Initialize θ = 0.
Randomly pick i = {1, ..., n}.
θ = (1 - ηλ) θ + η (y^{(i)} - θ · x^{(i)}) x^{(i)}.
```

**Hyperbolic Tangent Function**  $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = 1 - \frac{2}{e^{2z} + 1}$

**Deep Neural Networks** A deep (feedforward) neural network refers to a neural network that contains not only the input and output layers, but also hidden layers in between. Below is a deep feedforward neural network of 2 hidden layers, with each hidden layer consisting of 5 units:



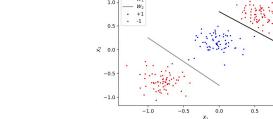
In the figure,

$$\tilde{W}_1 = \begin{pmatrix} W_{11} \\ W_{12} \end{pmatrix} \quad \text{and} \quad \tilde{W}_2 = \begin{pmatrix} W_{21} \\ W_{22} \end{pmatrix}.$$

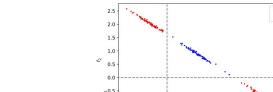
They map the input onto the  $f_1, f_2$  axes.

**Hidden Layer Representation**

• Hidden Layer Units



• Linear Activation



$f_1 = \tanh(z_1)$

$f_2 = \tanh(z_2)$

$f_3 = \tanh(z_3)$

$f_4 = \tanh(z_4)$

$f_5 = \tanh(z_5)$

$f = \tanh(z)$

$\mathcal{L}(f, f_1) = \text{Loss}(y, f_1) = \frac{1}{2} (y - f_1)^2$

$\mathcal{L}(f, f_2) = \text{Loss}(y, f_2) = \frac{1}{2} (y - f_2)^2$

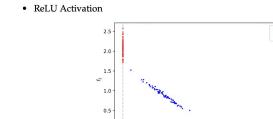
$\mathcal{L}(f, f_3) = \text{Loss}(y, f_3) = \frac{1}{2} (y - f_3)^2$

$\mathcal{L}(f, f_4) = \text{Loss}(y, f_4) = \frac{1}{2} (y - f_4)^2$

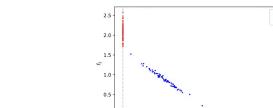
$\mathcal{L}(f, f_5) = \text{Loss}(y, f_5) = \frac{1}{2} (y - f_5)^2$

Note that the above derivation applies to tanh activation.

**Backpropagation** Consider the  $L$ -layer neural network below.



• tanh Activation



$f_1 = \tanh(z_1)$

$f_2 = \tanh(z_2)$

$f_3 = \tanh(z_3)$

$f = \tanh(z)$

**Summary**

- Units in neural networks are linear classifiers, just with different output non-linearities.
- The neurons in feedforward neural networks are arranged in layers.
- By learning the parameters associated with the hidden layer units, we learn how to represent examples (as hidden layer activations).
- The representations in neural networks are learned directly to facilitate the end-to-end task.
- A simple classifier (output unit) suffices to solve complex classification tasks if it operates on the hidden layer representations.

**Feedforward Neural Networks, Back Propagation, and Stochastic Gradient Descent (SGD)**

**Simple Example** This simple neural network is made up of  $L$  hidden layers, but each layer consists of only one unit, and each unit has activation function  $\text{f}$ .

$$x \rightarrow w_1 \text{---} z_1 \text{---} f_1 = \text{tanh}(z_1) \text{---} \dots \text{---} w_L \text{---} z_L \text{---} f_L = \text{tanh}(z_L)$$

$$\mathcal{L}(f_L, y_L) = \text{Loss}(y_L, f_L) = \frac{1}{2} (y_L - f_L)^2$$

For  $i = 2, \dots, L$ :  $z_i = f_{i-1} w_i$ , where  $f_{i-1} = f(z_{i-1})$ . Also,  $y$  is the true value and  $f_L$  is the output of the neural network. Then,  $\delta_L = \frac{\partial \mathcal{L}}{\partial f_L}$  denotes the error of neuron  $j$  in layer  $L$ . Then,  $\delta_i = \frac{\partial \mathcal{L}}{\partial f_i}$  denotes the full vector of errors associated with layer  $i$ .

If the activation function is  $f$  and the loss function we are minimizing is  $C$ , then the equations describing the network are

$$\delta_i^L = \frac{\partial \mathcal{L}}{\partial w_i} \odot \delta_{i+1}^L$$

$$\delta_i^L = [(w^{i+1})^T \delta^{i+1}]^T \odot f'(z^i)$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial C}{\partial f_i} \odot \delta_i^L$$

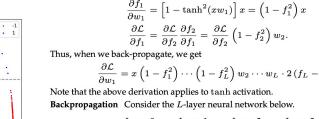
$$\frac{\partial \mathcal{L}}{\partial w_i} = (1 - f_i^2) \odot (1 - f_i^2)^T w_i \dots w_L \odot (f_L - y)$$

Thus, when we back-propagate, we get

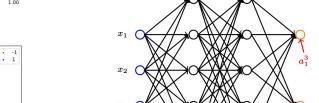
$$\frac{\partial \mathcal{L}}{\partial w_i} = (1 - f_i^2) \odot \dots \odot (1 - f_{i-1}^2) w_i \dots w_L \odot (f_L - y)$$

Note that the above derivation applies to tanh activation.

**Backpropagation** Consider the  $L$ -layer neural network below.



• ReLU Activation



$f_1 = \max\{0, z_1\}$

$f_2 = \max\{0, z_2\}$

$f_3 = \max\{0, z_3\}$

$f = \max\{0, z\}$

$\mathcal{L}(f, f_1) = \text{Loss}(y, f_1) = \max\{0, z_1 - y\}$

$\mathcal{L}(f, f_2) = \text{Loss}(y, f_2) = \max\{0, z_2 - y\}$

$\mathcal{L}(f, f_3) = \text{Loss}(y, f_3) = \max\{0, z_3 - y\}$

$\mathcal{L}(f, f) = \text{Loss}(y, f) = \max\{0, z - y\}$

Note that the above derivation applies to ReLU activation.

**Key Concepts**

- Encoding - e.g., mapping a sequence to a vector
- Decoding - e.g., mapping a vector to a sequence, e.g., a sequence

**Example: Encoding Sentences**

- Introduce adjustable **logo pieces** and optimize them for end-to-end performance.



context or state  $s_{t-1}$  →  $\oplus$  new information  $x_t$  → new context or state  $s_t$

$s_t = \tanh(W^{s,x} s_{t-1} + W^{s,x} x_t)$

$\oplus$  denotes element-wise addition.

**Nonlinear Classification**

**Feature Transformation**

$$x \rightarrow \phi(x)$$

$$\theta \cdot x \rightarrow \theta' \cdot \phi(x)$$

**Non-linear Classification**

$$h(x; \theta, \theta_0) = \text{sign}(\theta \cdot \phi(x) + \theta_0)$$

**Kernel Function** A kernel function is simply an inner product between two feature vectors.

Using kernels is advantageous when the inner products are faster to evaluate than using explicit vectors (e.g., when the vectors would be infinite dimensional).

Support Vector Machine finds the maximum margin linear separator by solving the quadratic program that corresponds to  $J(\theta, \theta_0)$ .

In the realizable case, if we disallow any margin violations, the quadratic program we have to solve is

Find  $\theta, \theta_0$  to minimize  $\frac{1}{2} \|\theta\|^2$  subject to

$$y^{(i)} (\theta \cdot x^{(i)} + \theta_0) \geq 1, \quad i = 1, \dots, n$$

**Stochastic Gradient Descent** Select  $i \in \{1, \dots, n\}$  at random

$$\theta \leftarrow \theta - \eta_i \nabla_\theta \left[ \text{Loss}_i(\theta \cdot x^{(i)} + \theta_0) + \frac{\lambda}{2} \|\theta\|^2 \right]$$

$\eta_i$  is the learning rate which can vary at every iteration.

**Support Vector Machine**

- Support Vector Machine finds the maximum margin linear separator by solving the quadratic program that corresponds to  $J(\theta, \theta_0)$ .

In the realizable case, if we disallow any margin violations, the quadratic program we have to solve is

Find  $\theta, \theta_0$  to minimize  $\frac{1}{2} \|\theta\|^2$  subject to

$$y^{(i)} (\theta \cdot x^{(i)} + \theta_0) \geq 1, \quad i = 1, \dots, n$$

**Nonlinear Classification**

**Decision Boundary** The decision boundary is the set of points  $x$  which satisfy

$$\theta \cdot x + \theta_0 = 0$$

**Margin** The margin is the set of points  $x$  which satisfy

$$\theta \cdot x + \theta_0 = \pm 1$$

**Hinge Loss** The hinge loss is

$$\text{Loss}_i(\theta$$

**First-order Markov Model** In a first-order Markov model (**bigram model**), the next symbol only depends on the previous one. Each symbol (except **(beg)**) in the sequence is predicted using the same condition probability table until an **(end)** symbol is seen. The reward probability associated to the sentence is

$$\prod_{i=1}^n P(w_i | w_{i-1}).$$

**Maximum Likelihood Estimation** The goal is to maximize the probability that the model can generate all the observed sentences (**corpus**)

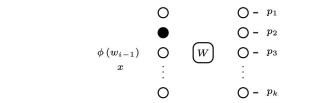
$$s \in S, s = \{w_1^*, w_2^*, \dots, w_n^*\}$$

$$\ell = \log \left\{ \prod_{s \in S} \prod_{i=1}^n P(w_i^* | w_{i-1}^*) \right\}$$

The maximum likelihood estimate is obtained as normalized counts of successive word occurrences (matching statistics)

$$\hat{P}(w' | w) = \frac{\text{count}(w', w)}{\sum_w \text{count}(w, w')}$$

**Feature-based Markov Model** We can also represent the Markov model as a feedforward neural network (very extensible). We define a one-hot vector,  $\phi(w_{i-1})$ , corresponding to the previous word. This will be an input to the feedforward neural network.



In the figure,

$$p_k = \mathbb{P}(w_k = k | w_{i-1})$$

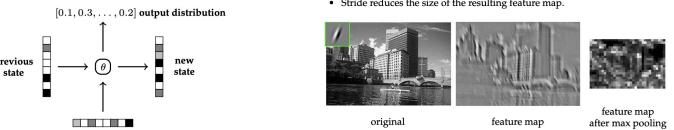
is the probability of the next word, given the previous word. The aggregate input to the  $k^{\text{th}}$  output unit is

$$z_k = \sum_k x_j W_{jk} + W_{0k}.$$

These input values are not probabilities. A typical transformation is the softmax transformation:

$$p_k = \frac{e^{x_k}}{\sum_j e^{x_j}}.$$

**RNNs for Sequences** Our RNN now also produces an output (e.g., a word) as well as update its state



Example of CNN From LeCun (2013)



## Clustering: Input

- Set of feature vectors  $S_n = \{x^{(i)} | i = 1, \dots, n\}$
- The number of clusters  $K$

## Clustering: Output

- A partition of indices  $\{1, \dots, n\}$  into  $K$  sets,  $C_1, \dots, C_K$
- Representatives in each of the  $K$  partition sets, given as  $z_1, \dots, z_K$

**Cost** We can calculate the total cost by summing the cost of each cluster:

$$\text{Cost}(C_1, \dots, C_K) = \sum_{j=1}^K \text{Cost}(C_j)$$

**Similarity Measure** We use the Euclidean distance between the elements of a cluster and its representative to calculate the cost for each cluster. Then, the total cost is

$$\text{Cost}(C_1, \dots, C_K, z_1, \dots, z_K) = \sum_{j=1}^K \sum_{i \in C_j} \|x^{(i)} - z_j\|^2.$$

## K-Means Algorithm

1. Randomly select  $z_1, \dots, z_K$ .
2. Iterate:

(a) Given  $z_1, \dots, z_K$ , assign each data point  $x^{(i)}$  to the closest  $z_j$  so that

$$\text{Cost}(z_1, \dots, z_K) = \sum_{i=1}^n \min_{j=1, \dots, K} \|x^{(i)} - z_j\|^2.$$

(b) Given  $C_1, \dots, C_K$ , find the best representatives  $z_1, \dots, z_K$ , i.e., find  $z_1, \dots, z_K$  such that

$$z_j = \arg \min_s \sum_{i \in C_j} \|x^{(i)} - z_j\|^2 = \frac{1}{|C_j|} \sum_{i \in C_j} x^{(i)}.$$

**K-Means Algorithm** The K-Means algorithm is a variation of the K-Means algorithm that addresses some of the K-Means algorithm's limitations.

1. Randomly select  $z_1, \dots, z_K \subseteq \{x_1, \dots, x_n\}$ .

2. Iterate:

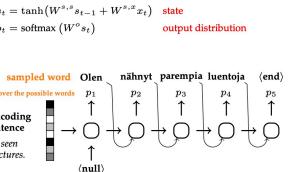
(a) Given  $z_1, \dots, z_K$ , assign each  $x^{(i)}$  to the closest  $z_j$  so that

$$\text{Cost}(z_1, \dots, z_K) = \sum_{i=1}^n \min_{j=1, \dots, K} \text{dist}(x^{(i)}, z_j)$$

(b) Given  $C_j \in \{C_1, \dots, C_K\}$ , find the best representative  $z_j \in \{x_1, \dots, x_n\}$  such that

$$\sum_{x^{(i)} \in C_j} \text{dist}(x^{(i)}, z_j)$$

is minimal.

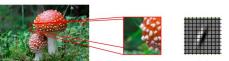


## Convolutional Neural Networks

**Problem** Image classification

- The presence of objects may vary in location across different images.

**Patch classifier/filter**



The patch classifier goes through the entire image. We can think of the weights as the image that the unit prefers to see.

**Convolution** The convolution is an operation between two functions  $f$  and  $g$ :

$$(f * g)(t) \equiv \int_{-\infty}^{+\infty} f(\tau)g(t - \tau)d\tau.$$

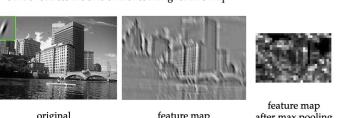
Intuitively, convolution **blends** the two functions  $f$  and  $g$  by expressing the amount of overlap of one function as it is shifted over another function.

**Discrete Convolution** For discrete functions, we can define the convolution as

$$(f * g)[n] \equiv \sum_{m=-\infty}^{m=\infty} f[m]g[n-m].$$

**Pooling** We wish to know whether a feature was there but not exactly where it was. Pooling (**Max**) Pooling region and stride may vary.

- Pooling induces translation invariance at the cost of spatial resolution.
- Stride reduces the size of the resulting feature map.



## Example: Image Quantization

## Unsupervised Learning

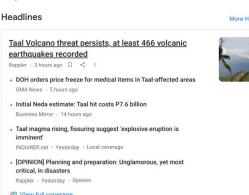
### Clustering

**Training set** We are provided a training set but with no labels

$$S_n = \{x^{(i)} | i = 1, \dots, n\}$$

and the goal is to find structure in the data.

**Example: Google News**



### Example: Image Quantization



**Partition** A partition of a set is a grouping of the set's elements into non-empty subsets, in such a way that every element is included in one and only one of the subsets. In other words,  $C_1, \dots, C_K$  is a partition of  $\{1, \dots, n\}$  if and only if

$$C_1 \cup \dots \cup C_K = \{1, \dots, n\} \quad \text{and}$$

$$C_i \cap C_j = \emptyset \quad \text{for any } i \neq j \in \{1, \dots, K\}.$$

## Generative Models

**Generative vs. Discriminative Models** Generative models work by explicitly modeling the probability distribution of each of the individual classes in the training data. Discriminative models learn explicit decision boundary classes.

**Simple Multinomial Generative Model** We assume that  $M$  is a fixed vocabulary  $W$  and we generate a document by sampling one word at a time from this vocabulary. Furthermore, all the words that are generated by  $M$  are independent of each other. We denote the probability that  $M$  generates certain word  $w \in W$  is

$$\mathbb{P}(w|\theta) = \theta_w, \quad \theta_w \geq 0, \quad \theta_w = 1.$$

Then, the probability of generating the document  $D$  is

$$\mathbb{P}(D|\theta) = \prod_{i=1}^n \theta_{w_i} = \prod_{w \in W} \text{count}(w).$$

**Maximum Likelihood Estimate** The log-likelihood for the model is

$$\ell = \log \mathbb{P}(D|\theta) = \sum_{w \in W} \text{count}(w) \log \theta_w$$

and the maximum likelihood estimate is

$$\hat{\theta}_w = \frac{\text{count}(w)}{\sum_{w' \in W} \text{count}(w')}.$$

**Prediction** Consider using a multinomial generative model  $M$  for the task of binary classification consisting of two classes: + (positive class) and - (negative class).

- $\theta^+$ : parameter for the positive class
- $\theta^-$ : parameter for the negative class

Suppose that we classify a new document  $D$  to belong to the positive class if and only if

$$\log \frac{\mathbb{P}(D|\theta^+)}{\mathbb{P}(D|\theta^-)} \geq 0.$$

The generative classifier is equivalent to a linear classifier:

$$\log \frac{\mathbb{P}(D|\theta^+)}{\mathbb{P}(D|\theta^-)} = \sum_{w \in W} \text{count}(w) \log \frac{\theta^+}{\theta^-} = \sum_{w \in W} \text{count}(w) \theta'_w.$$

**Prior, Posterior and Likelihood** In the above discussion, there is an assumption that the likelihood of being in one of the classes is the same. However, we may have some prior knowledge and we want to incorporate it into our model. The posterior distribution for the positive class is then

$$\mathbb{P}(y=+|D) = \frac{\mathbb{P}(D|\theta^+) \mathbb{P}(y=+)}{\mathbb{P}(D|\theta^+) + \mathbb{P}(D|\theta^-)}.$$

Prior knowledge  $\pi$  is the probability of being in the positive class. The posterior distribution is then

$$\mathbb{P}(y=+|D) = \frac{\pi \mathbb{P}(D|\theta^+) \mathbb{P}(y=+)}{\pi \mathbb{P}(D|\theta^+) + (1-\pi) \mathbb{P}(D|\theta^-)}.$$

**Gaussian Generative Models** The likelihood of  $\mathbf{x} \in \mathbb{R}^d$  being generated by a Gaussian with mean  $\mu$  and standard deviation  $\sigma$  is

$$f_X(\mathbf{x}|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mu\|^2\right).$$

### MLE for the Mean

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$

### MLE for the Variance

$$\hat{\sigma}^2 = \frac{1}{nd} \sum_{i=1}^n \|\mathbf{x}^{(i)} - \mu\|^2$$

**Markov Decision Processes**

**Definition** A Markov decision process (MDP) is defined by

- a set of states  $s \in S$  (may be observed or unobserved);
- a set of actions  $a \in A$ ;
- action-dependent transition probabilities  $T(s, a, s') = \mathbb{P}(s'|s, a)$  so that, for each state  $s$  and action  $a$ ,

$$\sum_{s' \in S} T(s, a, s') = 1$$

- reward functions  $R(s, a, s')$ , representing the reward for starting in state  $s$ , taking action  $a$  and ending up in state  $s'$  after one step. (The reward function may also depend only on  $s$ , or only on  $s$  and  $a$ .)

**Property** MDPs satisfy the **Markov property** in that the transition probabilities and rewards depend only on the current state and action, and remain unchanged regardless of the history (i.e., past states and actions) that leads to the current state.

**Utility Function** The main problem for MDPs is to optimize the agent's behavior. We first need to specify the criterion that we are trying to maximize in terms of accumulated reward. We define a utility function and maximize its expectation.

- **Finite horizon based utility**: the utility function is the sum of rewards after acting for a fixed number  $n$  of steps. When the rewards depend only on the states, the utility function is

$$U[s_0, s_1, \dots, s_n] = \sum_{i=0}^n R(s_i)$$

- **(Infinite horizon) discounted reward based utility**: In this setting, the reward one step into the future is discounted by a factor  $\gamma$ , the reward two steps ahead by  $\gamma^2$ , and so on. The goal is to continue acting (without an end) while maximizing the expected discounted reward. The discounting allows us to focus on near term rewards, and control the focus by changing  $\gamma$ . If the rewards depend only on the states, the utility function is

$$U[s_0, s_1, \dots] = \sum_{i=0}^{\infty} \gamma^i R(s_i)$$

**Optimal Policy** A policy is a function  $\pi : S \rightarrow A$  that assigns an action  $\pi(s)$  to any state  $s$ . Given an MDP and a utility function  $U[s_0, s_1, \dots, s_n]$ , our goal is to find an optimal policy function that maximizes the expectation of the utility. We denote the optimal policy by  $\pi^*$ .

### Bellman Equations

**Value Function** Denote by  $Q^*(s, a)$  the expected reward starting at  $s$ , taking action  $a$  and acting optimally. The value function  $V^*(s)$  is the expected reward starting at state  $s$  and acting optimally.

**The Bellman Equations** These equations connect the notion of the value of a state and the value of policy.

$$V^*(s) = \max_a Q^*(s, a) = Q^*(s, \pi^*(s))$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

We can define the  $V^*(s)$  recursively to get

$$V^*(s) = \max_a \left\{ \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')] \right\}$$

• Introduction to Machine Learning with Python (Müller and Guido)

• Machine Learning with Python – From Linear Models to Deep Learning [Lecture Slides] (<http://www.edx.org>)

• LaTeX File ([github.com/mynamiejanus/6.06MachineLearning](https://github.com/mynamiejanus/6.06MachineLearning))

Please share this cheatsheet with friends!

## Mixture Models; EM Algorithm

**Gaussian Mixture Models** Instead of just a single Gaussian, we have a mixture of Gaussian components. Assume that there are  $K$  Gaussians with known means and variances. Assume also that the mixture weights  $p_1, \dots, p_K$  are known. The likelihood for an observation  $x$  obtained from the model is

$$p(x|\theta) = \sum_{j=1}^K p_j \mathcal{N}(x; \mu^{(j)}, \sigma_j^2).$$

For the training set

$$S_n = \{x^{(i)} | i = 1, \dots, n\},$$

the likelihood is

$$\mathbb{P}(S_n|\theta) = \prod_{i=1}^n \sum_{j=1}^K p_j \mathcal{N}(x^{(i)}; \mu^{(j)}, \sigma_j^2).$$

**Observed Case** Consider the case of hard clustering, i.e., a point either belongs to a cluster or not. Let

$$\delta(j|i) = \begin{cases} 1, & x^{(i)} \text{ is assigned to } j \\ 0, & \text{otherwise.} \end{cases}$$

Also, let  $\bar{n}_j = \sum_{i=1}^n \delta(j|i)$  denote the number of points belonging to cluster  $j$ .

Maximizing the likelihood gives

$$\hat{p}_j = \frac{\bar{n}_j}{n}$$

$$\hat{\mu}^{(j)} = \frac{1}{\bar{n}_j} \sum_{i=1}^n \delta(j|i) x^{(i)}$$

$$\hat{\sigma}_j^2 = \frac{1}{\bar{n}_j d} \sum_{i=1}^n \delta(j|i) \|x^{(i)} - \mu^{(j)}\|^2.$$

**The EM Algorithm** Instead of hard clustering, the data can actually be generated from different clusters with different probabilities. We have soft clustering. We can maximize the likelihood through the EM algorithm.

Randomly initialize  $\theta: \mu^{(1)}, \dots, \mu^{(K)}, \sigma_1^2, \dots, \sigma_K^2, p_1, \dots, p_K$ .

1. **E-step:**

$$p(j|i) = \frac{p_j \mathcal{N}(x^{(i)}; \mu^{(j)}, \sigma_j^2 \mathbf{I})}{p(x|\theta)}$$

where  $p(x|\theta) = \sum_{j=1}^K p_j \mathcal{N}(x^{(i)}; \mu^{(j)}, \sigma_j^2 \mathbf{I})$

2. **M-step:**

$$\hat{n}_j = \sum_{i=1}^n p(j|i)$$

$$\hat{p}_j = \frac{\hat{n}_j}{n}$$

$$\hat{\mu}^{(j)} = \frac{1}{\hat{n}_j} \sum_{i=1}^n p(j|i) x^{(i)}$$

$$\hat{\sigma}_j^2 = \frac{1}{\hat{n}_j d} \sum_{i=1}^n p(j|i) \|x^{(i)} - \mu^{(j)}\|^2.$$

## Reinforcement Learning

**Objectives of RL** The goal of RL is to learn a good policy with no or limited supervision.

## Value Iteration Algorithm

**Definition** Value iteration is an iterative algorithm that computes the values of states indexed by  $k$ . Let  $V_k^*(s) = 0$  for  $s \in S$

$$V_k^*(s) \leftarrow V_{k-1}^*(s) + \gamma V_k^*(s)$$

$$V_k^*(s) \leftarrow \max_a \left\{ \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^*(s')] \right\}$$

$$V_k^*(s) \leftarrow \max_a Q_k^*(s, a)$$

**Convergence** This algorithm will converge as long as  $\gamma < 1$ .

## Q-Value Iteration

**Definition** We can directly operate at the level of Q-values. Q-value iteration is a reformulation of value iteration algorithm.

## Update Rule

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_a Q_k^*(s', a)]$$

## Reinforcement Learning

**MDP vs. RL** In MDPs, we are given 4 quantities  $(S, A, T, R)$ . In reinforcement learning, we are given only the states and actions  $(S, A)$ . In the real world, transitions and rewards might not be directly available and they need to be estimated.

**Estimation** Consider a random variable  $X$ . The goal is to estimate

$$E[f(X)] = \sum_x p(x) f(x).$$

We have access to  $K$  samples:  $x_1, x_2, \dots, x_K$ .

## Model-based Learning

$$\hat{p}(x_i) = \frac{1}{K} \text{count}(x_i)$$

$$E[f(X)] \approx \sum_{i=1}^K \hat{p}(x_i) f(x_i)$$

## Model-free Learning

$$E[f(X)] \approx \frac{1}{K} \sum_{i=1}^K f(x_i)$$