$$t > t_{n-k-1, \alpha/2} \text{; } n \text{ is samples}$$
$$k \text{ is predictors}$$
$$\alpha \text{ is confidence}$$

$$t = \frac{\beta_i - \beta^{null}(0)}{SE(\beta_i)}$$

when $\gamma \to \infty$, $t_r \to Z$, $Z$ is normal

$$CI = \beta \pm (t_{n-k-1, \alpha/2}) \cdot SE$$

null is $n-2$ DoF

$$F_{p, n-p-1} \underset{\alpha}{=} \frac{(TSS - RSS)/p}{RSS/(n-p-1)} \quad \begin{array}{l} p \text{ is} \\ \text{predictors} \end{array}$$

$$TSS - RSS = SSR \text{ (regression sum of squares)}$$
$$SSR = \Sigma (\hat{y}_i - \bar{y})^2$$

**Logistic Regression:**

$$P(x) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

False positive: Type 1
False negative: Type 2
Sensitivity: $TP/P$ (recall)
Specificity: $TN/N$
precision: $TP/TP+FP$
Neg predictive value: $TN/TN+FN$
$F_1$ score $= 2 \times \frac{precision \times recall}{precision + recall}$

$$F\beta = \frac{\frac{\beta^2 + 1}{1}}{\frac{\beta^2}{recall} + \frac{1}{precision}}$$

$$LOOCV = \frac{1}{n} \sum_{i=1}^{n} MSE_i$$

$$KFCV = \sum_{k=1}^{K} \frac{n_k}{n} MSE_k$$

**Logit:** "odds"

$$\log_{ln}\left(\frac{P(x)}{1 - P(x)}\right) = \beta_0 + \beta_1 X$$

$$Z = \frac{\hat{\beta}_i - \beta^{(null)}}{SE(\hat{\beta}_i)} \text{ standard normal}$$

Naive Bayes! LDA but w/
$$f_k(v) = \Pi_{j=1}^{P} f_{jk}(x_j)$$

**QDA**
$$d_k(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k) + \log \pi_k$$

**Lesson 4**
**Resampling**

Bootstrapping: $\bar{S}^* = \Sigma_b S(Z^{*b})/B$

$$\hat{Var}[s(z)] = \frac{1}{B-1} \sum_{b=1}^{B}(S(Z^{*b}) - \bar{S}^*)^2$$

$Z^*$ is a bootstrap dataset

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^{B}(\hat{\alpha}^{*r} - \bar{\hat{\alpha}}^*)^2}$$

$$\hat{Err}_{boot} = \frac{1}{B}\frac{1}{N} \sum_{b=1}^{B} \sum_{i=1}^{N} L(y_i, \hat{f}^{*b}(x_i))$$

True class1 | class2 (predicted c1 / c2)

$$RSE = \sqrt{\frac{1}{n-2} RSS} \text{ std of } \epsilon$$

$$TSS = \Sigma(y_i - \bar{y})^2$$
$$RSS = \Sigma(y_i - \hat{y}_i)^2$$

$$R^2 = \frac{TSS - RSS}{TSS}$$

$$r = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\Sigma(x_i - \bar{x})^2}\sqrt{\Sigma(y_i - \bar{y})^2}}$$

$$P_{Y|X}(y|x) = Pr(Y=y|X=x) = [p(x)]^y [1 - p(x)]^{1-y}, \quad y = 0, 1$$

**Bayes' Theorem**

$$Pr(Y=k|X=x) = \frac{Pr(X=x|Y=k) \cdot Pr(Y=k)}{Pr(X=x)}$$

**LDA** or $$pr(Y=k|X=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^{K} \pi_l f_l(x)}$$ ← density, prior

**LDA w/ $p=1$ Gaussian Density:**

$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2}\left(\frac{x - \mu_k}{\sigma_k}\right)^2}$$

or maximize

$$d_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

$$\hat{\mu}_k = \frac{1}{\pi_k} \sum_{i:y_i=k} x_i$$

$$\sigma^2 = \left(\frac{1}{n-k}\right) \sum_{k=1}^{K} \sum_{i:y_i=k}(x_i - \hat{\mu}_k)^2 = \sum_{k=1}^{K} \frac{n_k - 1}{n - k} \cdot \hat{\sigma}_k^2$$

$$\hat{Pr}(Y=k|X=x) = \frac{e^{d_k(x)}}{\sum_{l=1}^{K} e^{d_l(x)}}$$

**Linear Regression** $MSE = \frac{1}{n} \sum_{i=1}^{n}(y_i - \hat{f}(x_i))^2 = \frac{RSS}{n}$

$$Y = \beta_0 + \beta_1 X + \epsilon$$
$$\epsilon_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

$$SE(\beta_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$SE(\beta_0) = \sigma^2\left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right]$$

$$\sigma^2 = Var(\epsilon)$$

**Least squares!**

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} = \frac{\text{Sample covariance}}{\text{Sample variance}}$$

**Lesson 3**
**Classification**

$$p_k(x) = Pr(Y=k|X=x)$$

**Maximum Likelihood**

$$\ell(\beta_0, \beta) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0}(1 - p(x_i))$$

Choose $\beta_0 \cdot \beta_1$ to maximize

---

**Lesson 5: Model Selection + Regularization**

$$C_p = \frac{1}{n}\left(RSS + 2d\hat{\sigma}^2\right) \quad \begin{array}{l} \hat{\sigma}^2 \text{ is } \hat{var}(\epsilon) \\ d \text{ is \# parameters} \end{array}$$

$$\sigma(\overset{aka}{var(\epsilon)}) = RSE \quad \text{- hard to know}$$

$$\text{The Lasso: } RSS + \lambda \sum_{j=1}^{P} |\beta_j|$$

**Dimension Reduction**

$$Z_m = \sum_{j=1}^{P} \phi_{mj} X_j \text{; } Y_i = \theta_0 + \sum_{m=1}^{M} \theta_m Z_{im} + \epsilon \text{;}$$

$$\beta_j = \sum_{m=1}^{M} \theta_m \phi_{jm} \text{; } \theta \to \beta, Z \to X, Z \text{ via } \phi$$

$AIC = -2 \log L + 2 \cdot d$ (may like)

or

$$AIC = \frac{1}{n\hat{\sigma}^2}\left(RSS + 2d\hat{\sigma}^2\right)$$

$$BIC = \frac{1}{n}\left(RSS + \log(n) d\hat{\sigma}^2\right)$$

Elastic net for redundant variables. $L_1 + L_2$.

$$RSS + \lambda\left[\frac{1}{2}(1-\alpha)\|\beta\|_2^2 + \alpha\|\beta\|_1\right]$$

RR estimates should minimize: First standardize via:

$$\tilde{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{\sqrt{\frac{1}{n}\sum_{i=1}^{n}(x_{ij} - \bar{x}_j)^2}}$$

Choosing $\phi$ w/ PCR/PCA.
PLS considers Y too

**- Adjusted $R^2 = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)}$**

**Ridge Regression:**

$$RSS = \sum_{i=1}^{n}\left(y_i - \beta_0 - \sum_{j=1}^{P} \beta_j x_{ij}\right)^2$$

$$RSS + \lambda \sum_{j=1}^{P} \beta_j^2$$

# Lesson 6: Tree Based Methods

$$X_j < t_k$$

$$X_j < t_k \qquad X_j \geq t_k$$

· Find boxes that minimize
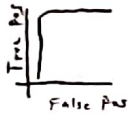$$RSS = \sum_{j=1}^{J} \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2 \quad (+ \alpha |T|)$$
← pruning

· cost-complexity pruning: grow a large tree, then prune it back, via cross-validation.

· Bagging: bootstrap (sample w/ replacement) our data repeatedly, getting B data sets. We average all predictions to obtain $\hat{f}_{bag}(x)$ (or maj vote).

· Boosting: Many trees, each fit to the residuals from previous trees. $\lambda$ is shrinkage parameter. B is # of trees - overfit!

1) set $\hat{f}(x) = 0$ + $r_i = y_i$ for all $i$
2) For $b = 1, 2 \ldots B$ : a) Fit tree $\hat{f}^b$ w/ d splits to training data $(X, r)$
   b) Update $\hat{f}$ by adding in shrunken version of new tree:
   $$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x) \qquad c) \text{ Update residuals } r_i - \lambda \hat{f}^b(x_i)$$
   d) Final Model: $\hat{f}(x) = \sum_{b=1}^{B} \lambda \hat{f}^b(x)$.

· Bayes Error Rate: $1 - E\left(\max_j Pr(Y=j \mid X)\right)$



· SMOTE: up + down sampling to account for imbalanced data.

· MSE = $Var\left(\hat{f}(x_0)\right) + Bias\left(f(x_0)\right) + Var(\varepsilon)$

· KNN: $p_j(x_0) = \frac{1}{k} \sum_{i \in N_0} I(y_i = j)$

---

Gini index: $G = \sum_{k=1}^{K} \hat{p}_{mk}(1 - \hat{p}_{mk})$    (classification) "purity"

A measure of total variance across k classes.
— or —

Cross-entropy: $D = -\sum_{k=1}^{K} \hat{p}_{mk} \log(\hat{p}_{mk})$.

· Random Forests: Decorrelated bagging. Only a random selection of m predictors is considered per-split. $m = \sqrt{p}$.

Stacking: Multiple models input to a meta model.
Adaptive meta learner: Implement different models on each section of the feature space, as is appropriate.

# Lesson 7: SVM's

A 2-D hyperplane: $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$

$X: \begin{cases} n \\ p \end{cases}$ is $X_{n,p}$.

Hence, $\beta_0 + \beta_1 x_1 + \ldots + \beta_p x_p \begin{cases} > 0, \text{ if } y_i = 1 \\ < 0 \text{ if } y_i = -1 \end{cases}$

or

$y_i(\beta_0 + \beta_1 x_1 + \ldots + \beta_p x_p) > 0$

Thus, we can classify a test observation based on the sign of $\beta_0 + \beta_1 x_1 + \ldots + \beta_p x_p$

- Magnitude suggests confidence

MMH is the solution to:

1) $\underset{\beta_0, \beta_1, \ldots \beta_p, M}{\text{maximize}} M$  | 2) Subject to: $\sum_{j=1}^{p} \beta_j^2 = 1$

3) $y_i(\beta_0 + \beta_1 x_{i_1} + \ldots + \beta_p x_{ip}) \geq M$ $\forall i=1,\ldots,n$

- M is the margin

**SVC**: Soft margin. Append above to:

$\geq M(1-\varepsilon_i)$; $\varepsilon_i \geq 0$; $\sum_{i=1}^{n} \varepsilon_i \leq C$ ← nonnegative tuning parameter (slack var)

**SVM**: Using polynomial functions of the predictors. Append above to.

$y_i\left(\beta_0 + \sum_{j=1}^{p} \beta_{j1} x_{ij} + \sum_{j=1}^{p} \beta_{j2} x_{ij}^2\right) \geq M(1-\varepsilon_i)$

- computationally infeasible, thus:

**Kernels**: $f(x) = \beta_0 + \sum_{i=1}^{n} \alpha_i \langle x, x_i\rangle$; where

$\langle x_i, x_{i'}\rangle = \sum_{j=1}^{p} x_{ij} x_{i'j}$ Alpha is nonzero only for the support vectors, $\sum_i \alpha_i = 0$.

Now, replace $\langle x, x_i\rangle$ with $K(x, x_i)$ (Kernel).

often $K(y_i, y_{i'}) = \left(1 + \sum_{j=1}^{p} x_{ij} x_{i'j}\right)^d$, but not always.

Thus now, $f(y) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$

- Multi-Class: One vs One (small K), One Vs All.
- Multi-Label (religions, politics, etc).
  Hamming Score: fraction of wrong labels to total # of labels.
- RBF: $K(x_i, x_{i'}) = \exp(-\gamma \sum(x_{ij} - x_{i'j})^2)$
- SVMs better than LR for $p \gg n$

---

# DSCI 552 Final

## Lesson 8: Unsupervised Learning

**K-means clustering**: Set K, then minimize $\left\{\sum_{k=1}^{k} W(C_k)\right\}$, "within cluster variation" which is average squared distance between all points in a cluster.

- Here's how: 1) Randomly assign each point.
2) Iterate: a) compute centroid, b) Assign to closest

Alternative: Metoid is centroid observation.

**Hierarchical Clustering**: Dendogram
1) Measure distance between all pairs. Each point is a cluster.
2) For $i = n, n-1, \ldots, 2$:
   a) Measure all pairwise inter-cluster dissimilarities & identify pair of clusters most similar. Fuse them.
   b) Repeat

Types of linkage: a) Complete: compute distance between all points in 2 clusters, record the largest
b) Single: record the smallest
c) Average: Average of all dissimilarities
d) centroid: Can result in undesirable inversions

- Correlation-based distance: Similar if features are correlated, even if euclid. fair.
- Between cluster variation: are groups spread apart? Overfits, find scree plot elbow
- Within cluster variation: optimize for small.
- Calinski-Harabasz index: Ideal local maximum for small W, large B.
  $CH(k) = \dfrac{B(k)/(k-1)}{W(k)/(n-k)}$ — maximize via K.
- Gap Statistic: How much W(k) drops @ each k.
  $G(k) = \log W_u(k) - \log W(k)$
- Silhouette analysis: $S_i = \dfrac{b_i - a_i}{\max(a_i, b_i)}$
  $a_i$ is average intracluster distance
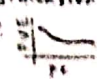  $b_i$ is average distance to nearest cluster

## principal component analysis

$Z_i = \phi_{11} X_1 + \phi_{21} X_2 + \ldots + \phi_{p1} X_p$ - find largest variance.
where $\sum_{j=1}^{p} \phi_{j1}^2 = 1$ is "normalized." (load gns ts)
- proportion of variance explained: PVE
  - What's lost by minimizing dimensions?
  - Visualize a scree plot:

- Fisher's LDA: supervised PCA.

---

## Lesson 9: Semi-Supervised Learning

- Transductive just creates labels, Inductive creates labels and a classifier.
- Self Training: inductive, classifier-based
  - A classifier is built on labeled data, used on unlabeled data, and the most confident are added in.
  - Refinement: reduce the weight of unlabeled data. "Yarowsky"
- Co-training: Build 2 classifiers on two different "views" of the data, if they agree, we add it.
- Cluster and label approach:
  - Cluster, classifier on labeled, assign unlabeled to each cluster label.
- Active Supervised Learning
  - Send intelligent queries about unlabeled data to an oracle who labels them. Update model.
  - What's "intelligent"?
    - Uncertain • Expected model change
    - Variance reduction
    - Query by committee: many models predict label, disagreements to oracle.

footer_navigationScanned with CamScannerfooter_navigation

# Lesson 10: Neural Networks

perceptron: $f(x) = \begin{cases} 1 & \beta^T x + \beta_0 \geq 0 \\ -1 & \beta^T y + \beta_0 < 0 \end{cases}$

Update rule: $\beta(i+1) = \beta(i) + 0.5\, e(i)\, x(i)$

$e(i) = y(i) - f(\beta^T(i) y(i) + \beta_0(i))$ - error

$\beta_0(i+1) = \beta_0(i) + 0.5\, e(i)$
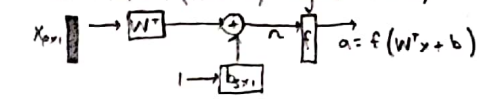
·Multi-class perceptron: (f is step-wise function)

· $W = \begin{bmatrix} w_1^T \\ \vdots & \vdots & S \\ w_i^T \end{bmatrix}$ ·Each class gets a hyperplane.

· $W(i+1) = W(i) + \alpha\, x(i)\, e^T(i)$

· $\alpha$ is learning rate or step size

· $n$: $W^T(i) x(i) + b(i)$  "net weight"

· $a$: $f(n) = f(W^T x + b) = \text{sign}(W^T x + b)$

$x_{p \times 1} \rightarrow \boxed{W^T} \rightarrow \oplus \rightarrow n \rightarrow \boxed{} \rightarrow a = f(W^T x + b)$
$1 \rightarrow \boxed{b_{y \times 1}}$

·Layers of perceptrons are for when the train set is not separable, new feature space.

· Sigmoid Function: $f(x) = \dfrac{1}{1+e^{-x}} = \dfrac{e^x}{e^x+1}$

· Tanh: $\dfrac{e^x - e^{-x}}{e^x + e^{-x}} = \dfrac{e^{2x}-1}{e^{2x}+1}$

· ReLu: ⟋  · M: # of layers

· J: objective function: something to be minimized by calculating weights. often "expected sum of square errors"

$E\left\{\sum_{i}(y_i - a_i^m)^2\right\} = E\{e^T e\}$

·Backpropagation Update Rule:

$S^m (\text{first}) = -2\dot{F}^{\prime(m)}(n^{(m)})\,|\,(y-a)$

$S^m (\text{subsequent}) = \dot{F}^{\prime(m)}(n^{(m)})\,W^{(m+1)} S^{(m+1)}$

where: $-\dot{F}^{\prime(m)}(n^{(m)}) = \dfrac{\delta f^m(n_j^m)}{\delta n_j^m}$

$\left(\begin{array}{l} \text{-note:} \\ \text{replace } -2(y-a) \text{ w/ } \nabla_a J \text{ w/ } J \neq e^T e \end{array}\right)$

· $W^{(m)}(k+1) = W^{(m)}(k) - \alpha a^{(m-1)} S^{(m)T}$
  (K from 1 to N)

· $b^{(m)T}(k+1) = b^{(m)T}(k) - \alpha S^{(m)T}$

## Regularization

In above, $W^{(m)}(k)(1-\eta\alpha)$ ← forgetting factor

· $\eta$ is decay rate

Empirical: Noisy input, noise to weights, rotate.

Softmax: Fixes sigmoid gradient problem

$p = e^a / (1+e^a)$ , also cross entropy.

X-entropy: $-y_2 \log a_2' - y_1 \log a_1'$ ← minimize

$p = a_2'$, $1-p = a_1'$, $y_1 = 1$, $y_2 = 0$

---

# Lesson 11: Hidden Markov Models

T: length of observation sequence

N: number of states in the model

M: number of observation symbols

Q: states of Markov process $(0 \rightarrow q_{N-1})$

V: set of possible observations $(0 \rightarrow M-1)$

A: State transition probabilities

B: observation probability matrix

$\pi$: initial state distribution

O: observation sequence

A is $N \times N$, B is $N \times M$

process: 1) Find p for each possible X
   - $O^N$ of them. For $\text{len}(O) = 4$,

eq: $\pi_{x_0} b_{x_0}(O_0) a_{x_0,x_1} b_{x_1}(O_1) a_{x_1,x_2} b_{x_2}(O_2) a_{x_2,x_3} b_{x_3}(O_3)$

The sum of all $O^N$ of them is the probability of O. Choose the highest for Dynamic Programming.

For Expectation Maximization, instead of highest total score, choose the likliest for each position.

## 3 types

1) Given $\lambda = (A, B, \pi) + O$, find $P(O|\lambda)$

2) Given $\lambda = (A, B, \pi) + O$, find optimal state sequence (hidden).

3) Given $O, N, M$, find $\lambda$ that maximizes probability of O.