

Final Exam

START TIME: 7:00 PM CST
STOP TIME: 10:00 PM CST

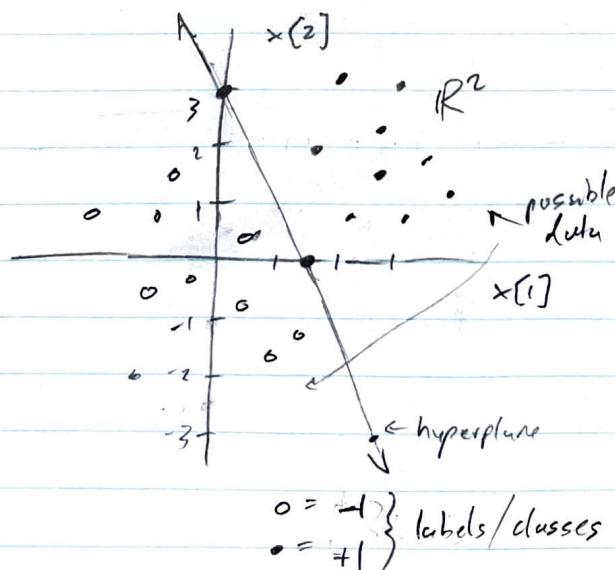
On my honor, I have neither given nor received
any unauthorized aid on this exam.

Robert Heeter

1. a. Regression refers to the set of problems/methods that involve designing a function that assigns continuous labels \hat{y}_i for input datapoints or vectors x_i . The labels \hat{y}_i are continuous / in \mathbb{R} . Classification refers to ~~a~~ a similar set of problems/methods, but where the labels \hat{y}_i are discrete (not continuous) and ~~fall~~ fall in a set $\hat{y}_i \in \{\dots\}$. While both regression & classification are ~~finite~~ types of supervised learning, classification methods require the ~~set of~~ finite set of labels to be predefined, whereas the ~~labels~~ labels in regression do not. Classification algorithms like logistic regression use special functions σ to map $\hat{y} \in \mathbb{R} \rightarrow \hat{y}_i \in \{\text{set of labels}\}$ $(0, 1)$

b. $w = \begin{bmatrix} 6 \\ 3 \end{bmatrix} \quad b = -9 \quad \hat{y}_i = \text{sign}(w^T x_i + b)$

hyperplane is \mathbb{R}^1 since w_1, b in \mathbb{R}^2



Intercepts: $x = [1.5, 0]^T, [0, 3]^T$

hyperplane:

$$[6, 3]^T \cdot x - 9 = 0$$

slope: -2

c. Next page

(i). c. PCA - principal components analysis - refers to the process of identifying the set of 1-dimensional subspaces of a data matrix that maximize the energy / spread / variance of the data which can then be used to transform the data matrix to lower dimensions that still retain the data's energy / spread / variance. The ~~present components are the first~~ ordered PCs are the leading ~~eigenvalues~~ eigenvectors of the data covariance matrix $X^T X$, where X is the data matrix, ~~and the~~ (in descending order) and the ~~eigenvalues~~ corresponding eigenvalues are the "energy" of that eigenvector (PC). In other words, to do PCA is to compute the PCs of the data matrix - ~~starting from~~ the eigenvectors of $X^T X$. This can be done with singular value decomposition $\rightarrow X = U \Sigma V^T$, where V 's columns are PCs. The PC decomposition is given as $P = X V = U \Sigma$. The ~~first~~ main PCs can be used to identify the most impactful features (dimensions in the original data matrix) since the size of the elements of ~~the~~ each eigenvector ~~is~~ corresponds to its ~~weight~~ - i.e. sorting the ~~the~~ eigenvector gives most to least ~~impactful~~ ^{feature's} features, some of which can be used to simplify the data matrix - reduce dimensions,

d. Multiclass support vector machines (SVM) can be performed in a "one vs - all" scheme. Here, a set of K 2-class SVMs are run ~~with~~ to ~~assign~~ assign the data, ^{vectors} to either one ~~of~~ of the K classes or the rest of the ~~the~~ classes. At the end, each data vector/point is labeled ~~with~~ according to the SVM that classified the point with the largest margin.

e. A nonlinear function ϕ can be applied to the columns of the data matrix X to better separate the two classes (i.e. so that they are now linearly separable) before using the linear classifier on the transformed data. \downarrow

$$x_i \rightarrow \phi(x_i) \Rightarrow \underline{x} \rightarrow \underline{\Phi(x)} \Rightarrow \text{use } \underline{\Phi} \text{ for linear classifier}$$

↑
column of \underline{X}

since data is now linearly separable.

$$f. \hat{y}_i = \sigma(\langle w, x_i \rangle + b)$$

In SVM, σ is typically simply the sign function - i.e. $\hat{y}_i = \text{sign}(\langle w, x_i \rangle + b)$ since the decision hyperplane $0 = \langle w, x_i \rangle + b$ separates data that is $\langle w, x_i \rangle + b > 0$ and < 0 . A single sum ~~is~~ decides between 2 ~~is~~ classes by defining the optimal ~~is~~/maximum margin hyperplane while sometimes accounting for slight errors (i.e. slack variables to define the tradeoff).

In logistic regression, σ is a scalar nonlinearity - i.e. a function that maps the result from the regression to ~~is~~ binary labels $[0, 1]$ $\in \mathbb{R}$

The sigmoid function $\sigma(u) = \frac{1}{1+e^{-u}}$ is an example.

$$y = \begin{cases} 1 & \text{if } \sigma(wx+b) \geq \frac{1}{2} \\ 0 & \text{if } \sigma(wx+b) < \frac{1}{2} \end{cases}$$

The hyperplane in logistic regression is defined from the orientation of the sigmoid function only from the ~~is~~ σ function nonlinear function, while the hyperplane in SVM is based on the geometrical and error calculation of the decision plane's margins and error.

g. Cross validation is important to determine how accurate a trained model fits the overall data. It works by dividing the data up into a set of groups, and then, for each group, training the model with the remaining (non-chosen group) data/groups and then testing that trained model on the chosen group. This is repeated for all of the groups such that all of the data is used for testing and the average accuracy of the model is averaged from each train/test process for each group. Other metrics to determine the model's "success" or performance can be used (not just accuracy). For example, with our final project for this class, cross-validation can be run with the set of "test" audio files - a fraction of the files can be saved for testing and the others used as training

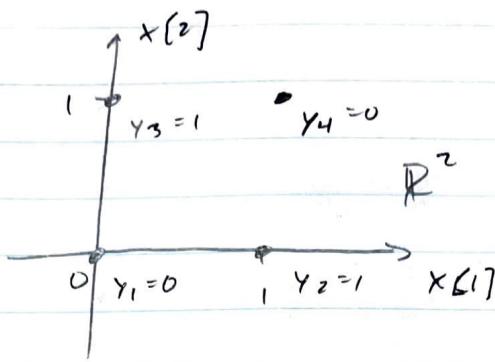
(i). (g), and the process can be repeated w/ another ~~set~~
now subset of the ~~the~~ "test" data, until all of the audio files
have been used for testing.

h. Kernel SVMs are advantageous in the fact that they ~~do~~ require
less training data than neural networks to produce a reasonable
decision boundary, and because they are ~~simpler~~ often
structurally simpler, it is ~~easier~~ to ~~be~~ computationally ~~more~~
less intensive to train a kernel SVM compared to a neural
network. ~~The~~ SVMs may be more useful for situations where there
is not as ~~much~~ much data available, since they still maximize the
decision boundary margins, while a neural network may have more
trouble in ^{i.e. builds a "buffer"} generalizing to a larger data set / real world situations.

Neural networks can be advantageous in situations where a lot of
data is ~~available and fast~~ and computational power is available
and accuracy is important since they can account for very subtle
features and patterns in the data using ~~a~~ gradual gradient descent

2. Next page

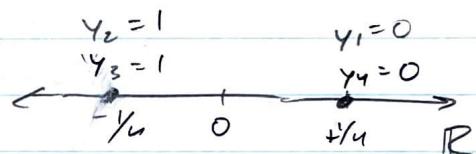
2. a.



$$\begin{cases} n = 4 \\ p = 2 \end{cases} \rightarrow 4 \text{ vectors} \quad \rightarrow 2 \text{ dimensions/features per vector}$$

- b. It is not possible for a single layer perceptron to correctly predict the label of all of the data points because the data is not linearly separable — i.e., a decision line boundary cannot be drawn that separates the $y=1$ and $y=0$ points, without error.

c. i	x_i	$\phi(x_i)$
1	$[0 \ 0]^\top$	$-1/2 \cdot -1/2 = +1/4$
2	$[1 \ 0]^\top$	$+1/2 \cdot -1/2 = -1/4$
3	$[0 \ 1]^\top$	$-1/2 \cdot +1/2 = -1/4$
4	$[1 \ 1]^\top$	$+1/2 \cdot +1/2 = +1/4$



$$\begin{cases} n = 4 \\ p = 1 \end{cases} \rightarrow 4 \text{ vectors} \quad \rightarrow 1 \text{ dimension/feature per vector}$$

d. $\phi(\underline{x}_i)$ is a nonlinear transform of \underline{x}_i .

To be a linear transformation: $\begin{cases} F(\underline{u} + \underline{v}) = F(\underline{u}) + F(\underline{v}) \\ \text{and } c F(\underline{v}) = F(c \underline{v}) \end{cases}$
 constant scalar

For example, $\underline{x}_1 = [0 \ 0]^\top \quad \underline{x}_1 + \underline{x}_2 = [1 \ 0]^\top$
 $\underline{x}_2 = [1 \ 0]^\top$

$$\phi(\underline{x}_1 + \underline{x}_2) = \phi([1 \ 0]^\top) = -1/4$$

$$\phi(\underline{x}_1) + \phi(\underline{x}_2) = \phi([0 \ 0]^\top) + \phi([1 \ 0]^\top) = +1/4 + (-1/4) = 0$$

$\therefore \phi(\underline{x}_1 + \underline{x}_2) \neq \phi(\underline{x}_1) + \phi(\underline{x}_2)$ and ϕ is not a linear transform of \underline{x} ;

$\therefore \phi$ is nonlinear transform of \underline{x} ;

(2). e. Yes, it is possible to use single layer perceptron to predict the label of the transformed data points.

$$\hat{z}_i = \sigma(v \cdot \phi(x_i) + c) \rightarrow \hat{z}_i \in \{0, 1\}$$

Choose $v = -4$, $c = 0$, $\sigma(u) = \begin{cases} 0 & \text{if } u < 0 \\ 1 & \text{if } u \geq 0 \end{cases}$ (step function)

i	x_i	$\phi(x_i)$	\hat{z}_i
1	$[0 0]^T$	$+4/4$	$-4(+4/4) = -1 < 0 \rightarrow 0$
2	$[1 0]^T$	$-4/4$	$-4(-4/4) = 1 \geq 0 \rightarrow 1$
3	$[0 1]^T$	$-4/4$	$-4(-4/4) = 1 \geq 0 \rightarrow 1$
4	$[1 1]^T$	$+4/4$	$-4(+4/4) = -1 < 0 \rightarrow 0$

$$v = -4$$

$$c = 0$$

$\sigma(u) = \text{step function (see above)}$

predictions
accurate

Another option w/ linear σ is:

$$(\hat{z}_i = \hat{y}_i)$$

$$v = -2$$

$$c = +1/2$$

$$\sigma(u) = u$$

$$\sigma(u) = \begin{cases} 0 & \text{if } u < 0 \\ 1 & \text{if } u \geq 0 \end{cases}$$

3. Next page

3. a. Measuring the error as $e = \frac{1}{2} (y_i - \hat{y}_i)^2$ is preferable to defining the error as $|y_i - \hat{y}_i|$ since the $\frac{1}{2}(y_i - \hat{y}_i)^2$ error is continuously differentiable, which is required for gradient descent to optimize the w and b that produce the lowest error. $|y_i - \hat{y}_i|$ is not differentiable at $\hat{y}_i = y_i$ whereas $\frac{1}{2}(y_i - \hat{y}_i)^2$ is.

b. The sigmoid function $\sigma(z) = \frac{1}{1+e^{-z}}$ may be preferred to the step function since it is continuously differentiable whereas the unit step function has a discontinuity at 0. This ~~sigmoid~~ is preferred for gradient descent, which requires that the loss function be continuously differentiable, since the ~~derivative~~ derivative is required for descent.

c. However, when $\sigma(z) = \frac{1}{1+e^{-z}}$, the predictions from $\hat{y} = \sigma(\langle x, w \rangle + b)$ are not binary, since $\sigma(z)$ does not return discrete values (i.e., not 1 or 0). The resulting predictions will never be truly binary (i.e., may take ~~values~~ values between 1 and 0).

d. $\hat{y} = \sigma(\langle x, w \rangle + b)$ with $\sigma(z) = \frac{1}{1+e^{-z}}$ is different from the logistic regression model since the logistic regression model restricts the predictions using inequalities:

$$\hat{y}_i = \begin{cases} 0 & \text{if } \sigma(\langle x_i, w \rangle + b) < \frac{1}{2} \\ 1 & \text{if } \sigma(\langle x_i, w \rangle + b) \geq \frac{1}{2} \end{cases}$$

Also, ~~the~~ logistic regression uses negative cross-entropy as an objective/loss function, rather than the ~~square~~ square of the error (half error squared) as is used here.

$$\rightarrow \text{i.e. Error} = \frac{1}{2} (y_i - \hat{y}_i)^2$$

e. Next page

$$(3). i.e. \nabla_{\underline{w}} \sigma(f(\underline{y})) = \underbrace{\sigma(f(\underline{y}))(1 - \sigma(f(\underline{y})))}_{\textcircled{1}} \nabla_{\underline{y}} f(\underline{y}) \quad \textcircled{2} \quad \textcircled{3}$$

$$\cancel{\nabla_{\underline{w}} l(\underline{w}, b)} = l(\underline{w}, b) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = \sigma(\underline{w}^T \underline{x}_i + b) = \sigma(\underline{w}^T \underline{x}_i + b)$$

$$l(\underline{w}, b) = \sum_{i=1}^n \frac{1}{2} (y_i - \sigma(\underline{w}^T \underline{x}_i + b))^2$$

$$\cancel{\nabla_{\underline{w}} l(\underline{w}, b)} = \cancel{\sum_{i=1}^n \frac{d}{d\underline{w}} (\sigma(\underline{w}^T \underline{x}_i + b))} \quad \downarrow$$

$$= \sum_{i=1}^n \left[\sigma(\underline{w}^T \underline{x}_i + b)(1 - \sigma(\underline{w}^T \underline{x}_i + b)) \right]$$

chain rule:

$$\begin{aligned} \nabla_{\underline{w}} l(\underline{w}, b) &= \sum_{i=1}^n \left[(y_i - \sigma(\underline{w}^T \underline{x}_i + b)) \left(0 - \underbrace{\sigma(\underline{w}^T \underline{x}_i + b)}_{\textcircled{1}} \underbrace{(1 - \sigma(\underline{w}^T \underline{x}_i + b))}_{\textcircled{2}} \underbrace{\frac{d}{d\underline{w}} (\underline{w}^T \underline{x}_i + b)}_{\textcircled{3}} \right) \right] \\ &= \sum_{i=1}^n (y_i - \sigma(\underline{w}^T \underline{x}_i + b)) \left[-\underline{x}_i \sigma(\underline{w}^T \underline{x}_i + b)(1 - \sigma(\underline{w}^T \underline{x}_i + b)) \right] \frac{d}{d\underline{w}} (\underline{w}^T \underline{x}_i + b) = \underline{x} \end{aligned}$$

$$\boxed{\nabla_{\underline{w}} l(\underline{w}, b) = -\sum_{i=1}^n \underline{x}_i (y_i - \sigma(\underline{w}^T \underline{x}_i + b)) \left[\sigma(\underline{w}^T \underline{x}_i + b)(1 - \sigma(\underline{w}^T \underline{x}_i + b)) \right]}$$

$$\nabla_b l(\underline{w}, b) = \sum_{i=1}^n \left[(y_i - \sigma(\underline{w}^T \underline{x}_i + b)) \left(0 - \underbrace{\sigma(\underline{w}^T \underline{x}_i + b)}_{\textcircled{1}} \underbrace{(1 - \sigma(\underline{w}^T \underline{x}_i + b))}_{\textcircled{2}} \underbrace{\frac{d}{db} (\underline{w}^T \underline{x}_i + b)}_{\textcircled{3}} \right) \right]$$

$$\boxed{\nabla_b l(\underline{w}, b) = -\sum_{i=1}^n (y_i - \sigma(\underline{w}^T \underline{x}_i + b)) \left[\sigma(\underline{w}^T \underline{x}_i + b)(1 - \sigma(\underline{w}^T \underline{x}_i + b)) \right]}$$

U.F. Next page

$$(3). f. \phi(x_i) \in \mathbb{R}^P \quad P \gg p$$

$$\phi(x_i) \in \mathbb{R}^P$$

A nonlinear transformation ϕ allows the linear classifier to potentially form a prediction for data that is nonlinearly separable, given an appropriate transformation/nonlinear func. ϕ . Using ϕ to transform the data $x_i \in \mathbb{R}^P$ and then calculate labels for the new matrix requires $O(P)$ which is more expensive than the $O(p)$ of using the original data. In some cases, the kernel trick can be used to avoid the costly computation of the P -dimensional transformed matrix and reduce the complexity back to $O(p)$, but then the kernel matrix and nonlinear transformation must satisfy Mercer's Theorem $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$. Not all ϕ can be used to make a kernel matrix and apply the kernel trick, in which case the more lengthy $O(P)$ calculation of the transformed data must be done.

$$g. \cancel{\alpha \in \mathbb{R}^n} \text{ such that } w = x^T \alpha \text{ where } X = [x_1 \dots x_n]^T$$

This assumption alone does not guarantee that the kernel trick can be applied, since there must exist an $\alpha \in \mathbb{R}^n$ such that $w = \phi^T \alpha$ where ϕ is the transformed data matrix X . This requires that the nonlinear function ϕ is appropriate to compute the transformation, and kernel matrix $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$

$$\hat{y}_i = \sigma(\langle x_i, w \rangle + b) \rightarrow \hat{y}_i = \sigma(\underbrace{\alpha \phi(x)}_{w}^T \underbrace{\phi(x_i)}_{\text{transformed } x} + b)$$

$$\hat{y}_i = \sigma(\underbrace{\alpha (\phi(x_i) \cdot \phi(x))}_{= K} + b)$$

\hat{y}_i can be found from

kernel matrix K , without

explicitly calculating transformed data matrix ϕ

$$\hat{y}_i = \sigma(\alpha K(x_i, x) + b)$$

Midterm Cheat Sheet

1. Linear algebra, optimization

Unsupervised: PCA, k-means clustering, hier. clust.

Supervised: lin reg, least squares, ridge/lasso reg.

2. Vector space V : a is scalar $\Rightarrow ax \in V, x+y \in V$

$$x, y \in V$$

Linear combination: $y = \sum_{m=1}^M a_m x_m = Xa$

vector \rightarrow i.e. mixing board

Euclidean 2 norm: $\|x\|_2 = \sqrt{\sum_{i=1}^n |x[i]|^2}$ \rightarrow energy $\|x\|_2^2$

Q -norm: $\|x\|_Q = (\sum_{i=1}^n |x[i]|^2)^{1/2}$

∞ -norm $\|x\|_\infty = \max_i |x[i]| \rightarrow$ peak value of x

Normalize vector $\|x\|_2 = 1 \rightarrow$ scale by $\frac{1}{\|x\|_2}$

Inner product: $x \cdot y = \langle x, y \rangle = y^T x = \sum_{i=1}^n x[i] y[i]$

Cosine similarity: $\cos(\theta)_{x,y} = \frac{x \cdot y}{\|x\|_2 \|y\|_2} \quad x \cdot x = \|x\|_2^2$

Orthogonality: $x \cdot y = 0 = \langle x, y \rangle$

3. Cauchy-Schwarz inequality: $0 \leq |x \cdot y| \leq \|x\|_2 \|y\|_2$

$x \neq y$ are: most diff or most similar
when orthogonal or when collinear

Basis: vector space V , lin. ind. & span V

dimension of V is p

$$b_k \cdot b_k = 0 \quad k \neq l$$

Orthogonal basis: $\{b_k\}_{k=1}^p \rightarrow$ elements orthogonal

Orthonormal basis \rightarrow orthogonal basis & normalized $\|b_k\|_2 = 1$

Inverse of matrix $AA^{-1} = A^{-1}A = I$

$\hookrightarrow B$ contains orthonormal basis, then $B^{-1} = B^T$

and B is orthogonal matrix

$$x = Ba = \sum_{k=1}^p a_k b_k \text{ synthesis; } a = B^T x, a_k = x \cdot b_k$$

Data matrix x $n \times p$ $[x]_{ij} \quad x_i: \text{datapoint (vector)}$

analysis

$x \in \mathbb{R}^p$ $p \times 1$ column vec.

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \begin{matrix} \text{vectors} \\ \text{rows} \end{matrix}$$

Correlation matrix $X^T X \rightarrow$ symmetric

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

Centering & normalizing: $\mu_i = \frac{1}{n} \sum_{l=1}^n x_l(i)$

$$x_l(i) \leftarrow \frac{x_l(i) - \mu_i}{\sigma_i} \quad \sigma_i = \sqrt{\frac{1}{n-1} \sum_{l=1}^n (x_l(i) - \mu_i)^2}$$

$$3. Av = \lambda v \quad \begin{matrix} \text{eigenvalues} \\ \uparrow \\ \text{must be square} \end{matrix}$$

$$A = V \Lambda V^{-1} \quad \begin{matrix} \text{real symmetric} \\ \text{matrix (correlation)} \\ \text{covariance} \end{matrix}$$

for example, have real-valued eigenvals & -vecs

Hermitian transpose A^H

$$A^H A = A A^H = I \rightarrow$$
 unitary matrix A (extension of orthogonal)

$$\text{SVD: } A = U \Sigma V^H$$

$$\begin{matrix} \text{np} & \text{np} & \text{np} \\ \uparrow & \uparrow & \uparrow \\ \text{unitary} & \text{diagonal} & \text{unitary} \\ \text{w/ eigenvectors} & \text{w/ singular values (non-reg)} & \text{w/ eigenvalues} \end{matrix} \quad \begin{matrix} \text{columns of } U, U^H \\ V, V^H \text{ are orthonormal bases} \end{matrix}$$

$$\text{Singular value } \sigma: Av = \sigma u, A^H u = \sigma u$$

$$\text{Frobenius norm } \|A\|_F^2 = \sum_{i=1}^p \sum_{j=1}^n |a_{i,j}|^2$$

Rank = # of nonzero singular vals.

4. PCA \rightarrow Data centered at 0 + normalized

$$\begin{matrix} x_i & \text{vector} \\ 0 & \text{origin} \\ u & \text{unit vector} \end{matrix} \quad \|u\|_2 = 1 \quad \varepsilon = \sum_{i=1}^n (x_i \cdot u)^2$$

$$\text{PCA} \rightarrow \max_{\|u\|_2=1} \sum_{i=1}^n (x_i \cdot u)^2 = \|x_c u\|^2$$

maximize energy/spread / variance of data after orthogonal projection onto

$$x^T X = V \Lambda V^T \rightarrow \max_{\|u\|_2=1} u^T V \Lambda V^T u \quad \text{(1st step)}$$

covariance matrix μ is first eigenvector v_1 of $x^T X$

$$\text{energy/eigenvector} \rightarrow \lambda_1 = \sum_{i=1}^n (x_i \cdot v_1)^2$$

2nd PC is orthogonal to first β is 2nd leading eigenvector

PCA via SVD: $x = U \Sigma V^T$

σ_k^2 is energy in PCs \rightarrow right singular vectors / columns of V are PCs

$$P = X V = U \Sigma$$

5. Objective (loss) func L , w_0 is minimizer/max...

Convex func: for any x_1, x_2 in F and $0 < t < 1$

$$f(tx_1 + (1-t)x_2) \leq t f(x_1) + (1-t) f(x_2)$$

Convex set: $t x_1 + (1-t)x_2 \in S$

$$\text{Local min } \|w - w^*\| \leq S \quad \text{constraint}$$

Lagrange multiplier $\lambda: \min_{w \in \mathbb{R}^p, \lambda \in \mathbb{R}} L(w) + \lambda g(w)$

$$\text{Gradient: } \nabla L(w) = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_p} \right]^T$$

$$= w^+ - \mu + \sum_{t=1}^T \nabla L_t(w)$$

↓
step size / learn rate

$$(5). \text{ Gradient descent: } w^{t+1} = w^+ - \mu^+ \nabla L(w^+)$$

↳ Issues w saddle and not smooth

Stochastic GD: $L(w) = \frac{1}{T} \sum_{t=1}^T L_t(w)$ ← sum of smaller obj. functions

↳ computationally less expensive; slow convergence

$$w^{t+1} = w^+ - \mu + \sum_{t=1}^T \nabla L(w)$$

↳ Epoch = all t used

set of gradients picked each iteration
(sequence or random)

6. Cluster → homogeneous groups; unsupervised

K-means: $\min_{C_1, \dots, C_K, w=1} \sum_{i=1}^n \|w(C_k)\|_2^2$ ← within cluster variation minimize

$$\rightarrow \|w(C_k)\|_2^2 = \sum_{i \in C_k} \|x_i - \bar{x}_k\|^2$$

↳ computation inefficient

finds local min

1) Random cluster assignment, 2) Centroid/ \bar{x}_k , 3) Assign to closest centroid

Different "distance" metrics
i.e. Euclidean $\sqrt{\sum (x_i - y_i)^2}$

1) $\|\cdot\|_2$ 2-norm

2) $\|\cdot\|_2^2$ squared 2

3) $\|\cdot\|_1$ 1-norm

4) $\|\cdot\|_\infty$ max norm
(max dist)

5) Cosine distance

Hierarchical clustering → tree-based taxonomy

↳ Agglomerative (bottom-up) ↳ Divisive (top-down)

↳ Similarity function → dendrogram tree

8. $x_i \rightarrow$ data point / vector $y_i \rightarrow$ label

$$y_i = w_1 + x_i w_2 + e_i \rightarrow \text{minimize error strength } \sum_{i=1}^n \|e_i\|^2$$

$$w_{LR} = \underbrace{\|y - xw\|^2}_{\text{objective fn}} \rightarrow L(w) = \|y - xw\|^2$$

$$\nabla L(w) = 0 \text{ @ minimizer } w^*$$

$$\nabla^2 L(w) = \left[\frac{\partial^2 L(w)}{\partial w_1}, \frac{\partial^2 L(w)}{\partial w_2} \right]^T$$

$$L(w) = \|y - xw\|^2 = (y - xw)^T (y - xw)$$

$$y^T \lambda w = w^T x^T y \rightarrow y^T y - y^T xw - w^T x^T y + w^T x^T xw \\ = y^T y - 2w^T x^T y + w^T x^T xw$$

$$\nabla_2 z^T a = a, \quad \nabla_2 z^T A_2 = A_2$$

$$\nabla^2 L(w) = -2x^T y + 2x^T xw = 0 \rightarrow x^T y = x^T xw$$

$$\text{Normal eqns: } x^T y = x^T xw^*$$

$$w^* = \underbrace{(x^T x)^{-1} x^T y}_{= x^T \text{ Moore Penrose}}$$

$$X = U\Sigma V^T, \quad X^T = V\Sigma^T U^T$$

↳ reciprocal of minor entries of Σ

$$\text{Quadratic} \quad X^T w + b$$

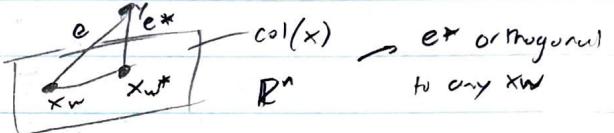
$$Y \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

$$8. w^* = (x^T x)^{-1} x^T y \quad \hat{y} = Xw^*$$



Orthogonality principle $x^T e^* = 0$

↳ optimal error vector e^* is orthogonal to each row of X^T and ∴ each column of X



Include column of 1s for y-intercept term w_0

Center data by subtracting mean of y and each column of x

$$I^T y = 0, I^T x = 0^T$$

→ w_0 is identifiable and $\propto w_0$'s scale normalized

$(x^T x)^{-1}$ does not exist → $x^T x$ is singular or ill-conditioned

1) $n \geq p$ (full, narrow x) but $\text{rank}(X) < p$

↳ 2+ columns of x (and $x^T x$) are lin. dependent

2) $n < p$ (short, wide x) → overfitting; more params w_j

↳ then data points (x_i, y_i) model fit to noise

LS Geometry axes are eigenvectors

w_1, w_2 eccentricity is from condition #

$$w_{LR} = w^*$$

Regularization w/ constraint: $\min_{w \in \mathbb{R}^p} \|w\|_1, \lambda \in \mathbb{R} \quad \|y - xw\|^2 + \lambda \|w\|_1^2$

Ridge $x^T y = (x^T x + \lambda I) w_{\text{ridge}} \rightarrow w_{\text{ridge}} = V(\Sigma + \lambda I)^{-1} U$
↳ always well-conditioned $\|w\|_2^2 = 0$

Lasso → 1 norm penalty $\|w\|_1 = c$

↳ Lasso is sparse, "variable selection"

Elastic net → combine 2 penalties

Need optimization alg. for Lasso

Final Cheat Sheet

11. SGD: $\mathcal{L}(w) = \sum_{m=1}^n \mathcal{L}_m(w)$

for ridge & lasso

select set of grads. (minibatch)

$$w^{t+1} = w^t - \mu \nabla \mathcal{L}_m(w^t)$$

① Ridge via GD:

$O(np)$

$$\stackrel{\text{min}}{w \in \mathbb{R}^p} \|y - Xw\|_2^2 + \lambda \|w\|_2^2 = \stackrel{\text{min}}{w \in \mathbb{R}^p} H(w)$$

$$\nabla H(w) = -X^T(y - Xw) + 2\lambda w$$

② Lasso via GD: (norm not differentiable)

$$\partial \mathcal{L} \approx 0$$

$$\stackrel{\text{min}}{w \in \mathbb{R}^p} \|y - Xw\|_2^2 + \lambda \|w\|_1 = \stackrel{\text{min}}{w \in \mathbb{R}^p} \mathcal{L}(w)$$

→ add condition for where $w_i = 0$

$O(np)$

12. Centered data; minimize error energy $\|y\|_2^2$

LMS

Time signals → window n samples $= \|y - Xw\|^2$

Wiener Filter

$$y = Xw + e \quad w = \text{impulse response}$$

Frobenius matrix, $Xw = \text{convolution}$

★ LMS algorithm: $w^{t+1} = w^t + \mu x_i (y_i - (x_i \cdot w^t))$

$O(p)$ ↳ System ID, linear prediction, L.P.C.

13. x_i → data point, feature, signal

Classification

y_i → label

$$y_i = f(x_i) \quad \begin{cases} \text{regression } y_i \in \mathbb{R} \\ \text{classification } y_i \in \{0, 1\} \\ \text{predictor} \end{cases} \quad \text{or layer set}$$

Decision boundaries

Error rates TPR, FPR, TNR, FNR



13
KNN

No optimization

1. Test data point x_{test}

3. $f(x_{test}) = \text{majority role of the } y_i$'s

2. Find K points closest

of K nearest neighbors

to x_{test} by some metric

• Good performance, handles many classes

• Computationally expensive $O(np)$ for $\frac{1}{2}$ test

→ Good for small data point set

15. $P(y_i = 1) = \sigma(w \cdot x + b)$

Logistic Reg

$$\sigma(u) = \frac{1}{1+e^{-u}}$$

$$y \begin{cases} 0 & \text{if } P(y=1) = \sigma(w \cdot x + b) < \frac{1}{2} \\ 1 & \text{if } P(y=1) = \sigma(w \cdot x + b) \geq \frac{1}{2} \end{cases}$$

↳ shifts sigmoid
sets classification threshold

multiple logistic regression

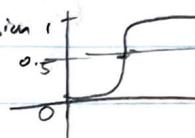
↳ w rules

hyperplane

level curves

of sigmoid surface in

RPM



↳ b → shifts sigmoid surface
in w direction

Objective fn: $\|y_i - \sigma(w \cdot x_i + b)\|_2^2 \rightarrow$ non convex

↳ use negative cross entropy due to sigmoid

$$\mathcal{L}(w, b) = \sum_{i=1}^n -y_i \log(\sigma(w \cdot x_i + b)) - (1-y_i) \log(1-\sigma(w \cdot x_i + b))$$

Derivative of sigmoid → convex!

$$\sigma'(z) = \frac{1}{1+e^{-z}} \rightarrow \sigma'(u) = \sigma(u)(1-\sigma(u))$$

Gradient step: $w^{t+1} = w^t + \mu \sum_{i=1}^n x_i (y_i - \sigma(w \cdot x_i + b))$

• Binary classifier

↳ one vs all classifier

for $K > 2$ classes

↳ classify each class vs rest of training

Multinomial Log. reg. (softmax)

↳ softmax nonlinearity: → convert K outputs to probabilities

$$\sum_{k=1}^K \exp(x_i \cdot w_k + b_k) \quad \left[\exp(x_i \cdot w_1 + b_1) \right]$$

↳ measure distances b/w histograms w/ negative cross entropy

objective fn: $\mathcal{L}(w, b, \dots, w_K, b_K)$

$$= -\sum_{i=1}^n \sum_{k=1}^K \frac{1}{\sum_j} \log \frac{\exp(x_i \cdot w_k + b_k)}{\exp(x_i \cdot w_j + b_j)}$$

Perceptron

$$\hat{y} = \sigma(x \cdot w + b) \quad \begin{matrix} \text{bias} \\ \text{applied component wise} \end{matrix}$$

scalar nonlinearities / activation function

$$\hat{y} = \sigma(w \cdot x + b) \quad \begin{matrix} \text{weight matrix} \\ \text{offsets vector} \end{matrix}$$

↳ sigmoid, hyperbolic tangent, ReLU

(15) MLP → multiple layers

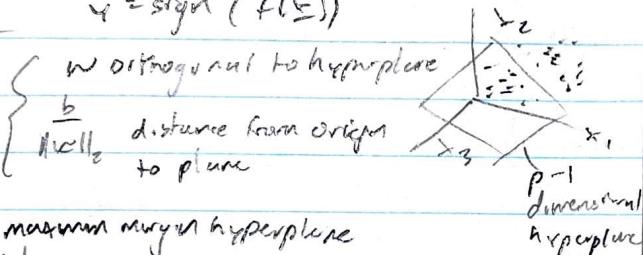
DNN → many params in $w_1 \dots w_L$ and $b_1 \dots b_L$
 ↳ constrain $w_1 \dots w_L$ to be constraint convolution
 matrices w/ short duration impulse response
 ↳ reduce # of params → CNN!
 ↳ Final activation σ^* → ReLU for regression
 σ^* → softmax for classif.

optimize
params
 w, b w/ SGD

backprop =
chain rule in calc

16. SVM → 2 classes

hyperplane $\mathcal{O} = f(\underline{x}) = \underline{w}^\top \underline{x} + b$
 $y = \text{sgn}(f(\underline{x}))$



#2 $K(K-1)/2$
 sum of w's
 choose each pair of
classes

choose class by most
sum selected
multiclass
SVM:
 #1 learn 1 class
vs all 2 classes

SVMs,
each class vs
rest of data
 given test data
point \underline{x} choose
sum w/ greatest
margin

large \rightarrow small
margin
but fewer
errors

Small \rightarrow opposite

driving hinge
loss
down drives
margin up

SVM
GD. O(np)

$\nabla_w J(w, b)$

$= \sum_{i=1}^n s_i^w(w) + 2\lambda w$

$\nabla_b J(w, b) = \sum_{i=1}^n s_i^b(b)$

column of X

replace $n \times p$ data matrix X w/ $n \times P$

data matrix $\Phi = \Phi(X)$

$x_i \rightarrow \phi(x_i)$

adds new term(s)

dim to data matrix

17. Invariant to nuisance transformations

Feature engineering → edge detection, cepstrum, SIFT

Feature learning → sparse coding

Deep learning → DNN

Cepstrum Features → cepstrum extracts elements

mel-freq. cepstrum

audio signals

$$F^{-1} [\log ((|Fx|^2)/2)]$$

Edge detection Fourier transform

↪ steep gradients

SIFT → scale invariant feature transform

dictionaries matrix $D \rightarrow x = \underline{Q}w$

sparse coding find w

$\min_w \|x - Dw\|^2 + \lambda \|w\|_1$

\uparrow $A^T F$ p x n transpose of data matrix

λ coeff. matrix

$\|\cdot\|_1$ 1-norm (sum of magnitude)

overcomplete dictionary Frobenius norm SSQures

bf convex

8. Ridge reg: $\min_w \|w\|^2 \|y - \underline{x}w\|^2 + \lambda \|w\|^2$

normal eqs $\underline{x}^\top y = (\underline{x}^\top \underline{x} + \lambda I) w$

primal sol'n $w = \underline{x}^\top \left(\frac{1}{\lambda} (y - \underline{x}w) \right) = \underline{x}^\top a$

dual sol'n $\min_{\alpha} \|\underline{y} - \underline{x}\alpha\|_2^2 + \lambda \|\underline{x}^\top \alpha\|_2^2$

$\underline{\alpha} = \underline{x}^\top T$ solve for α

α $n \times n$ dims (not $p \times p$)

$\hat{y} = \hat{f}(\underline{x}) = \underline{w}^\top \underline{x} = (\underline{x}^\top \underline{w})^\top \underline{x} = \underline{\alpha}^\top \underline{x}$

$\hat{y}_i = \hat{f}_i(\underline{x}) = \alpha_i^\top \underline{x}$

$\underline{\alpha} = \underline{x}^\top (\underline{x} \underline{x}^\top)^{-1} \underline{y}$

$K = \Phi \Phi^\top \rightarrow n \times n$

$[K]_{ij} = \phi(\underline{x}_i)^\top \phi(\underline{x}_j) = \phi(\underline{x}_i) \cdot \phi(\underline{x}_j)$

$\hat{y}_i = \underline{\alpha}^\top \underline{\Phi} \phi(\underline{x}) \rightarrow \hat{y}_i = \alpha_i \cdot \phi(\underline{x}_i) \cdot \phi(\underline{x})$

Mercer's Thm: class of funs ϕ such that

$\phi(\underline{x}_i) \cdot \phi(\underline{x}) = K(\underline{x}_i, \underline{x})$

Kernels: Polynomials $K(y, y) = (y \cdot y)^2$ mapping from

Gaussian \rightarrow radial basis fun don't need to work in D -dim

Sigmoid tanh space

19 \rightarrow Dropout

20. Feature selection \rightarrow statistical tests, correlation

Cross Validation \rightarrow divide matrix, PCA

overfitting as $P \rightarrow n$ coverage accuracies from each

test group

\rightarrow dropout, regularization, early stopping