

True/False

a) False, Naive Bayes assumes that all features are conditionally independent of the class. It's this assumption that limits its performance. It doesn't have to do with knowing the distribution.

b) True Linear regression is convex. But to have 1 min that is the global min, it must be strictly convex. This occurs when  $H$  is positive semidefinite

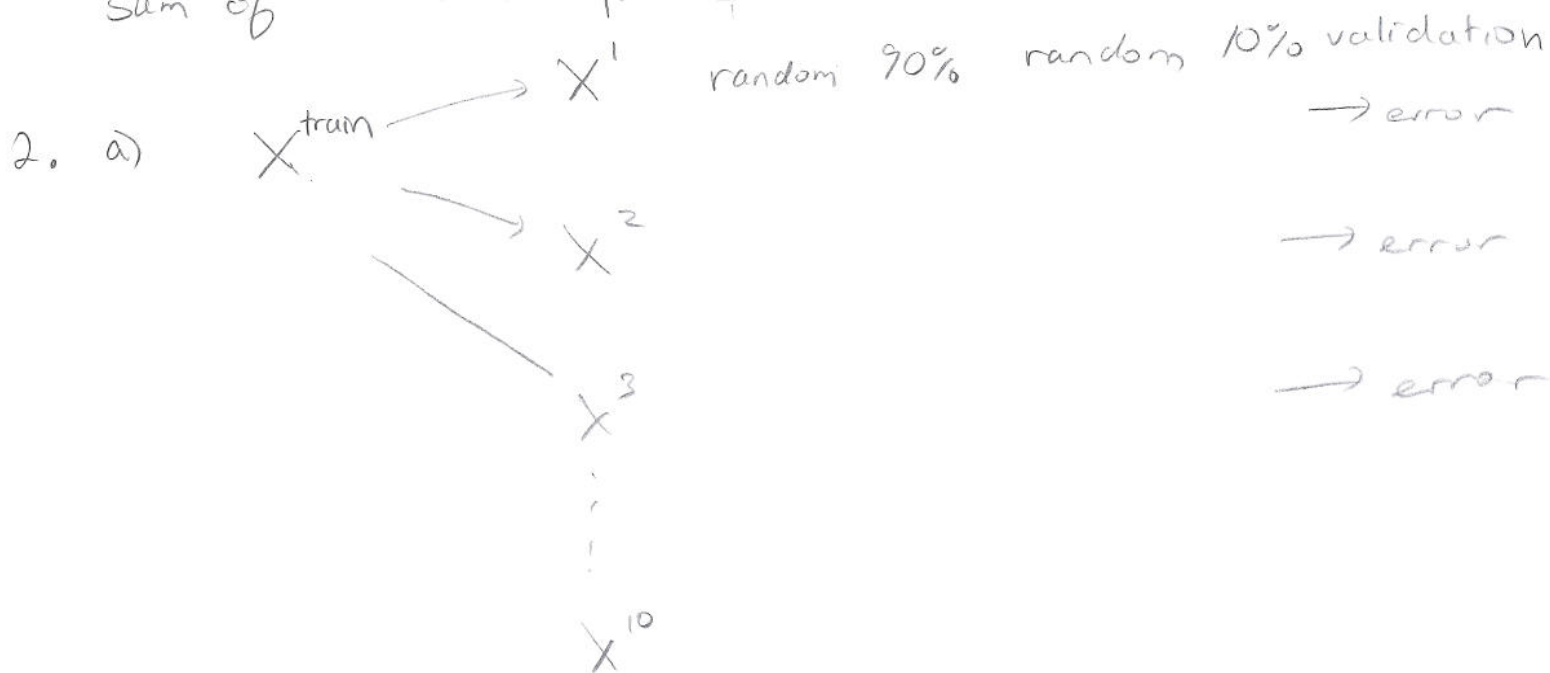
$$\frac{2}{2\beta} (-2(x^T y) + 2x^T x \beta)$$

$H = 2X^T X > 0$  only if  $X$  is full rank positive definite

- 1) convex  $\rightarrow$  multiple optima, same value
- 2) strictly convex  $\rightarrow$  1 optima, that is global

c) False, counter example is objective for  $l_1$ -regularized linear regression.  $l_1$  norm is not differentiable but objective is still convex. sq err +  $l_1$  norm convex, sum of convex functions is convex.

d) True, sq'd error +  $l_2$  norm are both convex. Sum of convex functions is convex.





get test error from test set

b) In the dual formation of the SVM, features appear only as dot products, which are represented compactly by kernels

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

c)  $\min x^2 + 1$

s.t.  $(x-2)(x-4) \leq 0$

$$L(x, \lambda) = x^2 + 1 + \lambda((x-2)(x-4)) \quad \text{where } \lambda \geq 0$$

$$= x^2 + 1 + \lambda(x^2 - 6x + 8)$$

$$= x^2 + 1 + \lambda x^2 - 6\lambda x + 8\lambda$$

$$= (\lambda + 1)x^2 - 6\lambda x + 8\lambda + 1$$

To get the dual we  $\max_{\lambda} (\min_x L(x, \lambda))$

$$\frac{\partial L(x, \lambda)}{\partial x} = 2(\lambda + 1)x - 6\lambda$$

$$0 = 2(\lambda + 1)x - 6\lambda$$

$$x = \frac{3\lambda}{\lambda + 1}$$

point that min L. via x.

b) In the dual of the SVM, features appear only at dot products, which are represented compactly by kernels

$$\max_{\alpha} \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j \underbrace{K(x_i, x_j)}$$

c)  $L(x, \lambda) = x^2 + 1 + \lambda((x-2)(x-4))$  where  $\lambda \geq 0$

$$= x^2 + 1 + \lambda(x^2 - 6x + 8)$$

$$= x^2 + 1 + \lambda x^2 - 6\lambda x + 8$$

$$= x^2 + 1 + \lambda x^2 - 6\lambda x + 8\lambda$$

$$= (\lambda + 1)x^2 - 6\lambda x + 8\lambda + 1$$

$\max_{\lambda} (\min_x L(x, \lambda))$

$$\frac{\partial L(x, \lambda)}{\partial x} = 2(\lambda + 1)x - 6\lambda$$

$$0 = 2(\lambda + 1)x - 6\lambda$$

$$x = \frac{6\lambda}{2(\lambda + 1)} = \frac{3\lambda}{\lambda + 1}$$

$$L(\lambda) = (\lambda + 1) \left( \frac{3\lambda}{\lambda + 1} \right)^2 - 6\lambda \left( \frac{3\lambda}{\lambda + 1} \right) + 8\lambda + 1$$

$$= \frac{9\lambda^2 (\cancel{\lambda + 1})}{(\lambda + 1)^2} - \frac{18\lambda^2}{\lambda + 1} + 8\lambda + 1$$

$$= \frac{9\lambda^2}{(\lambda + 1)} - \frac{18\lambda^2}{\lambda + 1} + 8\lambda + 1$$

$$= \frac{-9\lambda^2}{\lambda + 1} + 8\lambda + 1$$

$$L(\lambda) = (\lambda+1) \cdot \left(\frac{3\lambda}{\lambda+1}\right)^2 - 6\lambda \left(\frac{3\lambda}{\lambda+1}\right) + 8\lambda + 1$$

$$= \frac{9\lambda^2 (\cancel{\lambda+1})}{(\lambda+1)^{\cancel{2}}} - \frac{18\lambda^2}{\lambda+1} + 8\lambda + 1$$

$$= \frac{9\lambda^2}{(\lambda+1)} - \frac{18\lambda^2}{\lambda+1} + 8\lambda + 1$$

$$= \frac{-9\lambda^2}{\lambda+1} + 8\lambda + 1$$

$$\max \quad -\frac{9\lambda^2}{\lambda+1} + 8\lambda + 1$$

$$\text{s.t.} \quad \lambda \geq 0.$$

# Naive Bayes

	$P(x_a=1 y)$	$P(x_b=1 y)$	$x_c$	$x_d$	$x_e$	$x_f$	$x_g$
$y=1$	0	0	1/2	1	0	0	1/2
$y=-1$	1	1/2	0	0	1/2	1	0

$$P(y=1|z1) = 0$$

$$P(y=-1)$$

$$P(y=1|z1) = 1 \times \frac{1}{2} \times 1 \times \frac{1}{2} \times 1 \times 1 \times \frac{1}{2}$$

$$= \frac{1}{8}$$

$y=-1$

$$P(y=1|z2) = 1 \times 0$$

$$P(y=-1|z2) = 0$$

$x_1 =$	0	0	0	1	0	0	1	$y=1$
$x_2 =$	0	0	1	1	0	0	0	$y=1$
$x_3 =$	1	1	0	0	0	1	0	$y=-1$
$x_4 =$	1	0	0	0	1	1	0	$y=-1$
$x_5 =$	1	1	1	1	1	1	1	$y=1$
$x_6 =$	0	0	0	0	0	0	0	$y=1$
$x_7 =$	1	1	1	1	1	1	1	$y=-1$
$x_8 =$	0	0	0	0	0	0	0	$y=-1$

	$x_a$	$x_b$	$x_c$	$x_d$	$x_e$	$x_f$	$x_g$
$y=1$	1/4	1/4	2/4	3/4	1/4	1/4	2/4
$y=-1$	3/4	2/4	1/4	1/4	2/4	3/4	1/4

Name Bayes

$$P(y=1|z1) = \left(\frac{1}{4} \cdot \frac{3}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{1}{4} \cdot \frac{2}{4}\right) \cdot \frac{1}{2} = \frac{36}{2^{15}} \quad P(y=1)$$

$$P(y=-1|z1) = \left(\frac{3}{4} \cdot \frac{2}{4} \cdot \frac{3}{4} \cdot \frac{3}{4} \cdot \frac{2}{4} \cdot \frac{3}{4} \cdot \frac{3}{4}\right) \cdot \frac{1}{2} = \frac{3^5 \cdot 2^2}{2^{15}}$$

$$y = -1$$

$$P(y=1|z2) = \left(\frac{3}{4} \cdot \frac{1}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{1}{4} \cdot \frac{2}{4}\right) \frac{1}{2} = \frac{36}{2^{15}}$$

$$P(y=1|z2) = \left(\frac{1}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{2}{4} \cdot \frac{3}{4} \cdot \frac{1}{4}\right) \frac{1}{2} = \frac{36}{2^{15}}$$

he

Perceptron

Start

$$w = (0, 0, 0, 0, 0, 0, 0)$$

---

$$y_1 [\langle w, x_1 \rangle] = 0 \leq 0? \text{ Yes}$$

$$w = w + y_1 x_1 = (0, 0, 0, 1, 0, 0, 1)$$

---

$$y_2 [\langle w, x_2 \rangle] = 0 \leq 0? \text{ Yes}$$

$$w = w + y_2 x_2$$

$$= (0, 0, 0, 1, 0, 0, 1) + (-1, -1, 0, 0, 0, 0, -1, 0)$$

$$= (-1, -1, 0, 1, 0, -1, 1)$$

---

$$y_3 [\langle w, x_3 \rangle] = 1 \leq 0? \text{ No}$$

no update

---

$$y_4 [\langle w, x_4 \rangle] = (-1 - 1) - 1 = 2 \leq 0? \text{ No}$$

No update

---

$$y_5 [\langle w, x_5 \rangle] = (1 + 1) - 1 = 2 \leq 0? \text{ No}$$

No update

$$w = (-1, -1, 0, 1, 0, -1, 1)$$

# SVMs

1. Write down the problem

$$\min_{w, b} \frac{1}{2} \|w\|^2 + c \sum \xi_i$$

$$\text{s.t. } y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

for large values of  $c$ , penalizing shrinking the margin heavily,

that is penalizing misclassified points

∴ decision boundary will separate data perfectly if possible

2.  $c=0$ , not penalizing misclassified points at all. penalty is low, so we can misclassify a few while maximizing the margin btw most of the points

3. Warning was don't trust any specific data point too much, so we prefer  $c \approx 0$

4. Correctly classified by the original classifier, will not be a support vector

5.  $c$  is large, adding a point that is incorrectly classified by the original boundary would force the boundary to move.



3. [2 points] Which of the two cases above would you expect to work better in the classification task? Why?
  
4. [3 points] Draw a data point which will not change the decision boundary learned for very large values of  $C$ . Justify your answer.
  
5. [3 points] Draw a data point which will significantly change the decision boundary learned for very large values of  $C$ . Justify your answer.

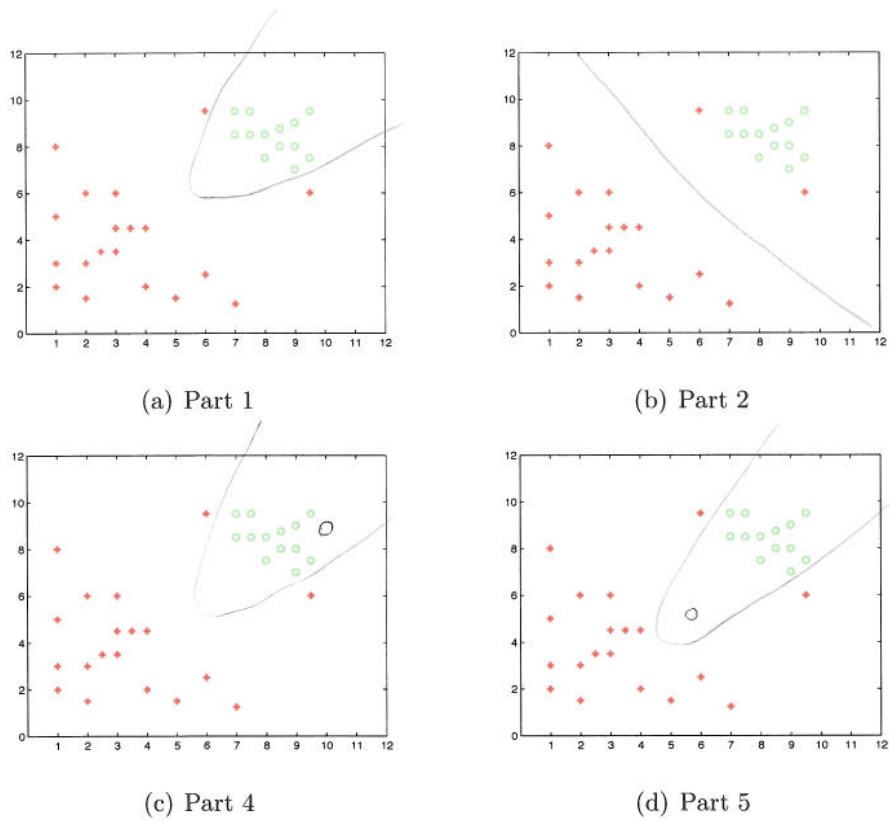


Figure 2: Draw your solutions for Problem 2 here.

# Conditional Independence, MLE/MP, Prob.

1. Use Chain Rule

$$P(A_n \dots A_1) = P(A_n | A_{n-1} \dots A_1) P(A_{n-1} \dots A_1)$$

$$\begin{aligned} P(x, y | z) &= P(x | y, z) P(y | z) \\ &= P(x | z) P(y | z) \end{aligned}$$

2.

$$L = \prod_{i=1}^n \frac{\theta^{y_i} e^{-\theta}}{y_i!}$$

$$\log L = \log \prod_{i=1}^n \frac{\theta^{y_i} e^{-\theta}}{y_i!}$$

$$= \sum_{i=1}^n (y_i \log \theta - \theta - \log y_i!)$$

$$= \log \theta \sum_{i=1}^n y_i - n\theta - \sum_{i=1}^n \log y_i!$$

$$\begin{aligned} 3. \quad P(\text{correct} | \text{answer}) &= 1 & P(\text{answer}) &= p \\ P(\text{correct} | \text{guess}) &= \frac{1}{m} & P(\text{guesses}) &= 1 - p \end{aligned}$$

$$P(\text{answer} | \text{correct}) = \frac{P(\text{correct} | \text{answer}) P(\text{answer})}{P(\text{correct})}$$

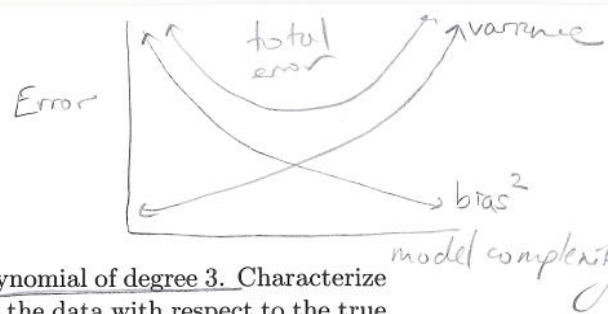
$$= \frac{1 \cdot p}{1 \cdot p + \frac{1}{m} (1-p)}$$

$$= \frac{p}{p + \frac{1}{m} (1-p)}$$

$$= \frac{P(\text{correct} | \text{answer}) P(\text{answer})}{P(\text{correct} | \text{answer}) P(\text{answer}) + P(\text{correct} | \text{guess}) (P(\text{guess}))}$$

$$\text{Err} = \text{Bias}^2 = (E[\hat{f}(x)] - f(x))^2$$

$$\text{Var} = [\hat{f}(x) - E[\hat{f}(x)]]^2$$



#### 4 Bias-Variance Decomposition (12 pts)

1. (6 pts) Suppose you have regression data generated by a polynomial of degree 3. Characterize the bias-variance of the estimates of the following models on the data with respect to the true model by circling the appropriate entry.

skip

	Bias	Variance
Linear regression	low/high	low/high
Polynomial regression with degree 3	low/high	low/high
Polynomial regression with degree 10	low/high	low/high

2. Let  $Y = f(X) + \epsilon$ , where  $\epsilon$  has mean zero and variance  $\sigma_\epsilon^2$ . In  $k$ -nearest neighbor (kNN) regression, the prediction of  $Y$  at point  $x_0$  is given by the average of the values  $Y$  at the  $k$  neighbors closest to  $x_0$ .

- (a) (2 pts) Denote the  $\ell$ -nearest neighbor to  $x_0$  by  $x_{(\ell)}$  and its corresponding  $Y$  value by  $y_{(\ell)}$ . Write the prediction  $\hat{f}(x_0)$  of the kNN regression for  $x_0$  in terms of  $y_{(\ell)}, 1 \leq \ell \leq k$ .

$$\hat{f}(x_0) = \frac{1}{k} \sum_{\ell=1}^k y_{(\ell)}$$

- (b) (2 pts) What is the behavior of the bias as  $k$  increases?

decreases

Solutions online are wrong

- (c) (2 pts) What is the behavior of the variance as  $k$  increases?

increases

## 5 Support Vector Machine (12 pts)

Consider a supervised learning problem in which the training examples are points in 2-dimensional space. The positive examples are  $(1, 1)$  and  $(-1, -1)$ . The negative examples are  $(1, -1)$  and  $(-1, 1)$ .

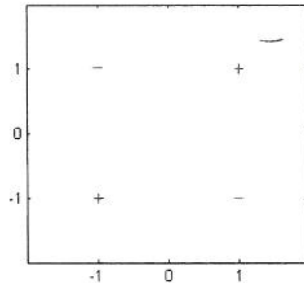
- (1 pts) Are the positive examples linearly separable from the negative examples in the original space?

No

- (4 pts) Consider the feature transformation  $\phi(x) = [1, x_1, x_2, x_1x_2]$ , where  $x_1$  and  $x_2$  are, respectively, the first and second coordinates of a generic example  $x$ . The prediction function is  $y(x) = w^T * \phi(x)$  in this feature space. Give the coefficients,  $w$ , of a maximum-margin decision surface separating the positive examples from the negative examples. (You should be able to do this by inspection, without any significant computation.)

$$w = (0, 0, 0, 1)^T$$

- (3 pts) Add one training example to the graph so the total five examples can no longer be linearly separated in the feature space  $\phi(x)$  defined in problem 5.2.



- (4 pts) What kernel  $K(x, x')$  does this feature transformation  $\phi$  correspond to?

$$\begin{aligned} & [1, x_1, x_2, x_1x_2] \\ & [1, x_1', x_2', x_1'x_2'] \end{aligned} \quad \phi(x) \cdot \phi(x')$$

7

$$1 + x_1x_1' + x_2x_2' + x_1x_1'x_2x_2'$$

## 6 Generative vs. Discriminative Classifier (20 pts)

Consider the binary classification problem where class label  $Y \in \{0, 1\}$  and each training example  $X$  has 2 binary attributes  $X_1, X_2 \in \{0, 1\}$ .

In this problem, we will always assume  $X_1$  and  $X_2$  are conditional independent given  $Y$ , that the class priors are  $P(Y = 0) = P(Y = 1) = 0.5$ , and that the conditional probabilities are as follows:

$P(X_1 Y)$	$X_1 = 0$	$X_1 = 1$	$P(X_2 Y)$	$X_2 = 0$	$X_2 = 1$
$Y = 0$	0.7	0.3	$Y = 0$	0.9	0.1
$Y = 1$	0.2	0.8	$Y = 1$	0.5	0.5

The expected error rate is the probability that a classifier provides an incorrect prediction for an observation: if  $Y$  is the true label, let  $\hat{Y}(X_1, X_2)$  be the predicted class label, then the expected error rate is

$$P_{\mathcal{D}}(Y = 1 - \hat{Y}(X_1, X_2)) = \sum_{X_1=0}^1 \sum_{X_2=0}^1 P_{\mathcal{D}}(X_1, X_2, Y = 1 - \hat{Y}(X_1, X_2)).$$

Note that we use the subscript  $\mathcal{D}$  to emphasize that the probabilities are computed under the true distribution of the data.

\*You don't need to show all the derivation for your answers in this problem.

- (4 pts) Write down the naïve Bayes prediction for all the 4 possible configurations of  $X_1, X_2$ . The following table would help you to complete this problem.

$X_1$	$X_2$	$P(X_1, X_2, Y = 0)$	$P(X_1, X_2, Y = 1)$	$\hat{Y}(X_1, X_2)$
0	0	$0.7 \cdot 0.9 \cdot 0.5$	$0.2 \cdot 0.5 \cdot 0.5$	0
0	1	$0.7 \cdot 0.1 \cdot 0.5$	$0.2 \cdot 0.5 \cdot 0.5$	1
1	0	$0.3 \cdot 0.9 \cdot 0.5$	$0.8 \cdot 0.5 \cdot 0.5$	1
1	1	$0.3 \cdot 0.1 \cdot 0.5$	$0.8 \cdot 0.5 \cdot 0.5$	1

- (4 pts) Compute the expected error rate of this naïve Bayes classifier which predicts  $Y$  given both of the attributes  $\{X_1, X_2\}$ . Assume that the classifier is learned with infinite training data.

skip

(a) By regularizing  $w_2$  [3 pts] Increases

By regularizing  $w_2$ , the boundary can rely less and less on  $x_2$  and thus boundary becomes more vertical.

(b) By regularizing  $w_1$  [3 pts] Same.

By regularizing  $w_1$ , the boundary can rely less and less on  $x_1$ , and thus boundary becomes more horizontal. That's ok b/c training data can be separated by horizontal linear separator.

(c) By regularizing  $w_0$  [3 pts]

Increase, When we regularize  $w_0$ , then the boundary will eventually go through the origin.

Best we can get is one error

2. If we change the form of regularization to L1-norm (absolute value) and regularize  $w_1$  and  $w_2$  only (but not  $w_0$ ), we get the following penalized log-likelihood

$$\max \sum_{i=1}^n \log P(y_i | x_i, w_0, w_1, w_2) - C(|w_1| + |w_2|).$$

Consider again the problem in Figure 1 and the same linear logistic regression model  $P(y = 1 | \vec{x}, \vec{w}) = g(w_0 + w_1 x_1 + w_2 x_2)$ .

- (a) [3 pts] As we increase the regularization parameter  $C$  which of the following scenarios do you expect to observe? (Choose only one) Briefly explain your choice:
- First  $w_1$  will become 0, then  $w_2$ .
  - First  $w_2$  will become 0, then  $w_1$ .
  - $w_1$  and  $w_2$  will become zero simultaneously.
  - None of the weights will become exactly zero, only smaller as  $C$  increases.

1) we can classify with zero error on  $x_2$  alone so  $w_1$  goes to zero. Note absolute value reg. ensures it goes exactly to zero

2) as  $C$  increases, we pay higher and higher cost for  $w_2$  so it eventually goes to zero

skip

(b) [3 pts] For very large  $C$ , with the same L1-norm regularization for  $w_1$  and  $w_2$  as above, which value(s) do you expect  $w_0$  to take? Explain briefly. (Note that the number of points from each class is the same.) (You can give a range of values for  $w_0$  if you deem necessary).

skip

(c) [3 pts] Assume that we obtain more data points from the '+' class that corresponds to  $y=1$  so that the class labels become unbalanced. Again for very large  $C$ , with the same L1-norm regularization for  $w_1$  and  $w_2$  as above, which value(s) do you expect  $w_0$  to take? Explain briefly. (You can give a range of values for  $w_0$  if you deem necessary).

# Kernel Regression

1.

$$\frac{2 J(\beta)}{2\beta} = -2 A^T W (Y - A\beta)$$

$$\frac{2}{2s} (x - A_s)^T W (x - A_s)$$

$$= -2 A^T W (x - A_s)$$

$$0 = -2 A^T W Y + 2 A^T W A \beta$$

$$A^T W Y = A^T W A \beta$$

$$\hat{\beta} = (A^T W A)^{-1} A^T W Y$$

2. when  $A^T W A$  is invertible, full rank

3. Gradient Descent

$$x^{(k+1)} = x^{(k)} - t_{k+1} \nabla f(x^{(k)})$$

$$\beta^{(t+1)} = \beta^{(t)} - \underset{\substack{\uparrow \\ \text{step size } 2\beta}}{\alpha} \cdot \frac{2 J(\beta)}$$

$$= \beta^{(t)} - \alpha (-2 A^T W (Y - A\beta))$$

$$= \beta^{(t)} - \alpha A^T W (A\beta - Y)$$

4.

$$Y = \beta_1 + \beta_2 X + e$$

$$e \sim N(0, \sigma_i^2)$$

$$\sigma_i^2 \propto \frac{1}{W_i(x)}$$



5. 1 advantage  $\rightarrow$  no strict assumptions on the form of the underlying distribution or regression function

1 disadvantage  $\rightarrow$  computationally expensive  
large # of training examples