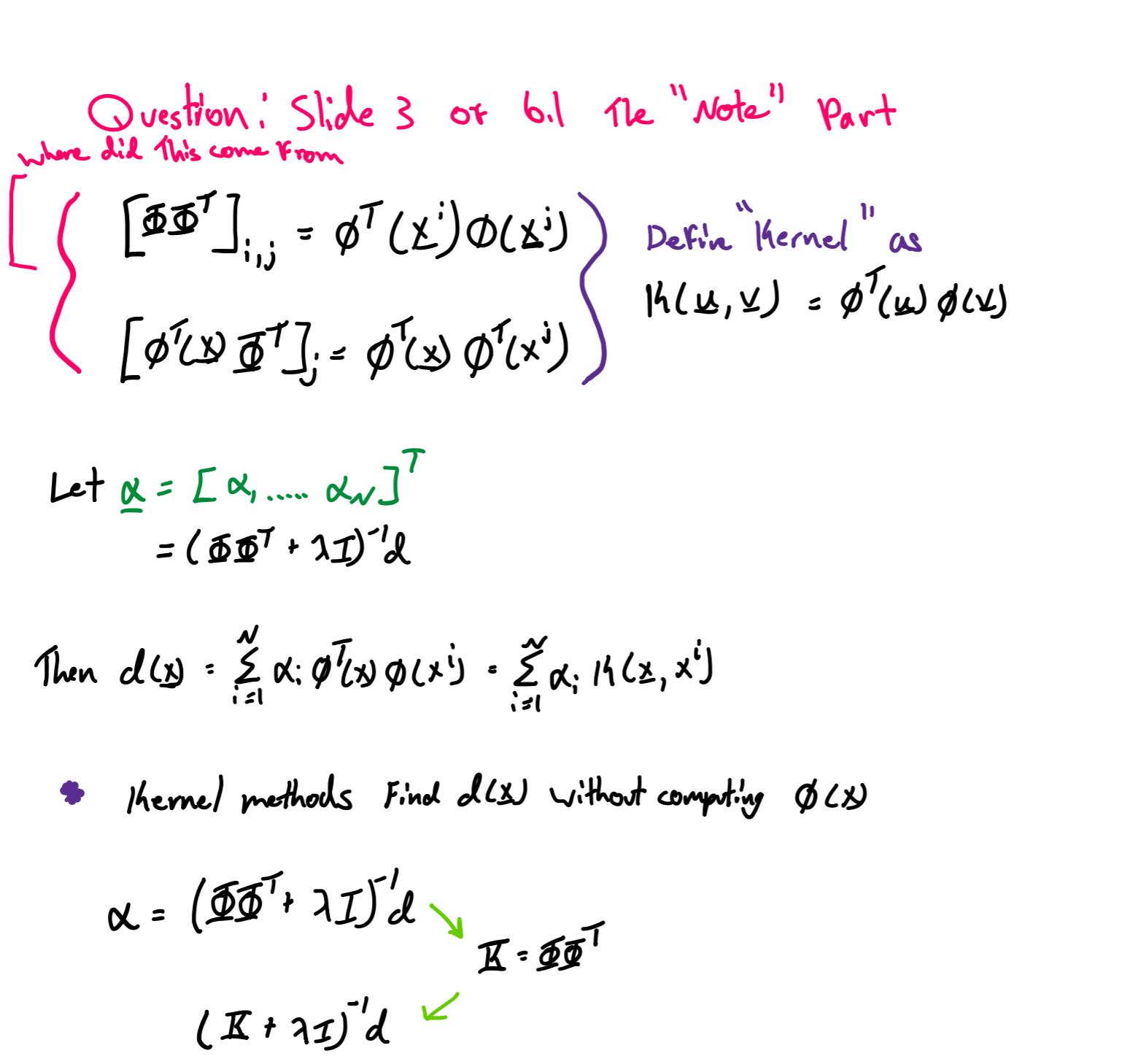


b.1 Kernel Regression

goals

- why use higher dim feature spaces
- reformulate regression in terms of kernels
- look at some popular kernels
- cautions and considerations

we let $\underline{x} = [x_1, x_2, x_3, \dots, x_n]^T \in R^n$ The variable $d(\underline{x}) = \phi^T(\underline{x})\underline{w}$, where $\phi(\underline{x}) \in R^{p > n}$
regression or modelexample) $\underline{x} = [x_1, x_2]^T$, $\phi(\underline{x}) = [x_1^2, x_2^2, \sqrt{x_1}x_2, x_1, x_2, 1]$ we find \underline{w} using training data. Through the process below

- Ridge regression: $\min_{\underline{w}} \frac{1}{2} (\underline{d} - \underline{w})^T (\underline{d} - \underline{w}) + \lambda \|\underline{w}\|_2^2$
- where $\underline{d} = [d_1, d_2, \dots, d_m]^T$
- so our sol is $\underline{w} = (\phi^T \phi + \lambda I)^{-1} \underline{d}$

now, $d(\underline{x}) = \phi^T(\underline{x})\underline{w} = \phi^T(\underline{x})(\phi^T \phi + \lambda I)^{-1} \underline{d}$
which we can be written as a weighted sum of "kernels"
Regressors

$$\text{Matrix Identity: } (\phi^T \phi + \lambda I)^{-1} \phi^T = \phi^T (\phi \phi^T + \lambda I)^{-1}$$

$$(\phi^T \phi + \lambda I)(\phi^T \phi + \lambda I)^{-1} \phi^T = \phi^T (\phi \phi^T + \lambda I)^{-1} (\phi^T + \lambda I)$$

thus, $d(\underline{x}) = \phi^T(\underline{x}) \underbrace{\phi^T (\phi \phi^T + \lambda I)^{-1}}_{N \times N} \phi$

Question: Slide 3 or b6. The "Note" Part
Where did this come from?

$$\begin{cases} [\phi \phi^T]_{i,j} = \phi^T(x^i) \phi(x^j) \\ [\phi^T \phi]_{i,j} = \phi^T(\underline{x}) \phi(\underline{x}) \end{cases}$$

Define "Kernel" as $K(\underline{x}, \underline{y}) = \phi^T(\underline{x}) \phi(\underline{y})$

Let $\underline{x} = [\alpha_1, \dots, \alpha_N]^T$
 $= (\phi \phi^T + \lambda I)^{-1} \underline{d}$

Then $d(\underline{x}) = \sum_{i=1}^N \alpha_i \phi^T(\underline{x}) \phi(x^i) = \sum_{i=1}^N \alpha_i K(\underline{x}, x^i)$

• Kernel methods find $d(\underline{x})$ without computing $\phi(\underline{x})$

$$\underline{x} = (\phi \phi^T + \lambda I)^{-1} \underline{d} \quad \underline{x} \cdot \underline{x}^T$$

$$[\underline{x} \cdot \underline{x}^T]_{i,j} = \phi^T(x^i) \phi(x^j)$$

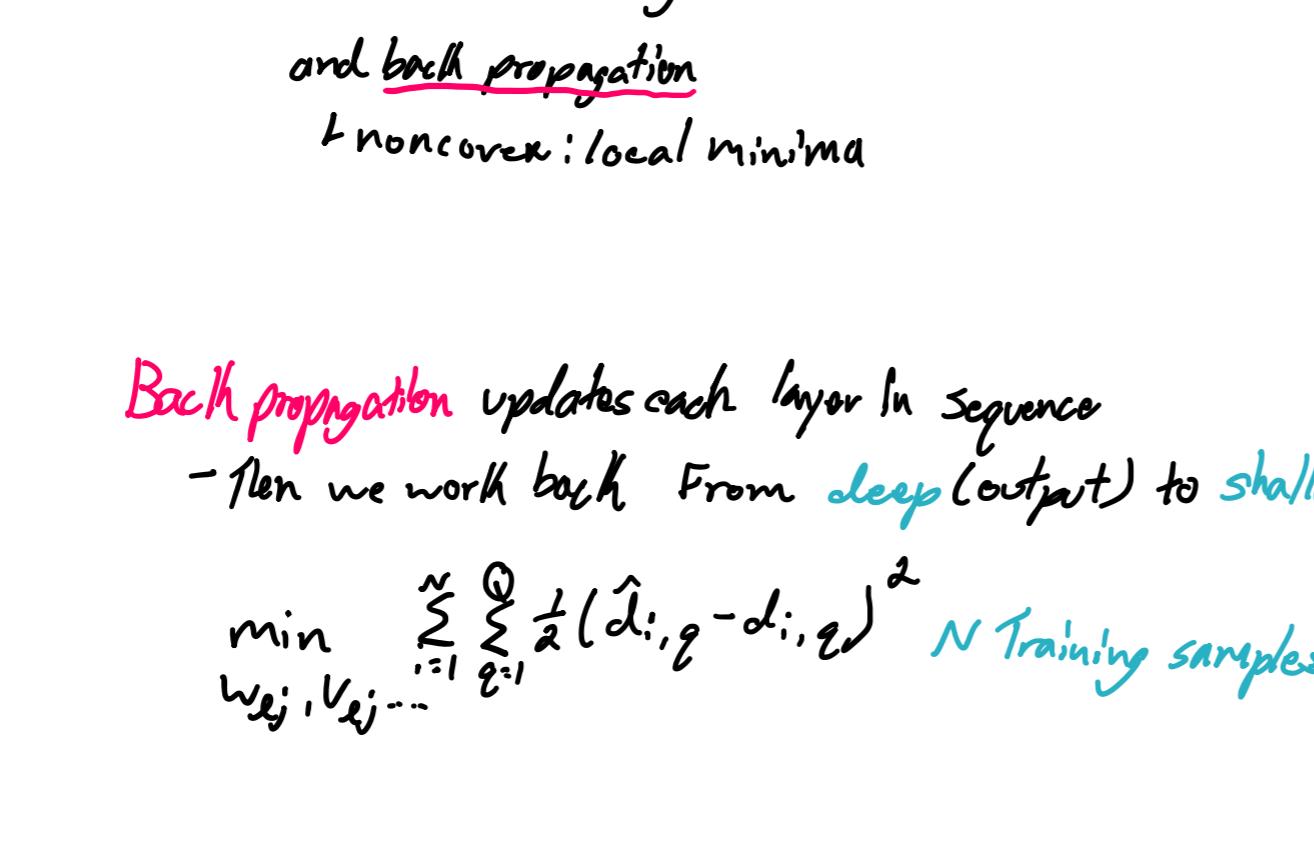
not really sure how but, \underline{x} can be computed efficiently!

- ex) Monomials of degree 2 Powers add to 2
 \downarrow
 $d(\underline{x}) \rightarrow x_1^2, x_1 x_2, \dots, x_2^2, x_3^2$

$$K(\underline{x}, \underline{y}) = \phi^T(\underline{x}) \phi(\underline{y}) = (\underline{x}^T \underline{y})^2$$

$$P = \frac{(e+m-1)!}{e!(m-1)!} \text{ terms}$$

Suppose
 $m=10, e=5 \rightarrow P \sim 2 \times 10^9$ • computing $O(10)$ vs $O(m)$
 $m=100, e=5 \rightarrow P \sim 10^8$ • memory $O(NP)$ vs $O(N^2)$

look at some popular kernels which one we use depends on $\underline{x}, \underline{y}$ monomials of degree 0: $K(\underline{x}, \underline{y}) = (\underline{x}^T \underline{y})^0$ monomials of degree 1: $K(\underline{x}, \underline{y}) = (\underline{x}^T \underline{y})^1$ Polynomials up to degree e : $K(\underline{x}, \underline{y}) = (\underline{x}^T \underline{y})^e$ Gaussian/radial kernel: $K(\underline{x}, \underline{y}) = \exp\left\{-\frac{\|\underline{x}-\underline{y}\|^2}{2\sigma^2}\right\}$ 

cautions and considerations

$$d(\underline{x}) = \phi^T(\underline{x}) \underline{w} \quad d(\underline{x}) = \sum_{i=1}^n \alpha_i K(\underline{x}, x^i)$$

store and compute $\alpha_i (N \times 1)$ vs $\underline{w} (1 \times N)$

binary class sign{d(x)} better boundaries

avoid over fitting with cross-validation CV

b.2 Kernel Based SVM

Goals:

- reformulate linear max margin classifier in terms of support vectors

- derive Kernel version of hinge loss with ridge regression

- summarize features of support vector machines

- linear max margin classifier
- If we have data points that are linearly separable we can derive a classifier that splits the data based on the max difference between data points. The points that are used to find max differences are support vectors.

$$w = [w^T, b]^T \text{ depends only on the support vectors}$$

$$\text{Recall Kernel regression: } d(\underline{x}) = \phi^T(\underline{x}) \underline{w} = \sum_{i=1}^n \alpha_i K(\underline{x}, x^i)$$

What happens when we apply kernel regression when we don't have a higher dimensional space? Well, Let $\phi(\underline{x}) = \underline{x} \Rightarrow K(\underline{x}, \underline{x}') = \underline{x}^T \underline{x}'$

$$\text{Then } d(\underline{x}) = \sum_{i=1}^n \alpha_i K(\underline{x}, x^i) = \sum_{i=1}^n \alpha_i x^T x^i$$

we can see the explicit dependence on \underline{x} in the training data. Also note that All $\alpha_i \neq 0$ except for support vectors

we can use kernels for non-linear decision boundaries

high-dim Feature space: $\underline{x} \rightarrow \phi(\underline{x})$

$$\text{ex) } \phi(\underline{x}) = [x_1^2, x_1, \dots, x_n^2, \dots, x_n, 1]$$

$$d(\underline{x}) = \text{sign}(\phi(\underline{x}))$$

Hinge loss with ridge regression

$$\min_{\underline{w}} \sum_{i=1}^n (1 - d(x^i))_+ + \lambda \|\underline{w}\|_2^2$$

Claim! $\underline{w} = \sum_{i=1}^n \phi(x^i) \alpha_i$

Now we rewrite the Hinge loss with ridge regression

Problem by replacing \underline{w} and replacing $\phi(x^i) \phi(x^j)$ with $K(x^i, x^j)$

$$\min_{\underline{w}} \sum_{i=1}^n (1 - d(x^i))_+ + \lambda \sum_{i=1}^n \alpha_i \sum_{j=1}^n \alpha_j K(x^i, x^j)$$

$$\min_{\underline{w}} \sum_{i=1}^n (1 - d(x^i))_+ + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x^i, x^j)$$

$$\text{SVM } \min_{\underline{w}} \sum_{i=1}^n (1 - \sum_{j=1}^n \alpha_j \alpha_i K(x^i, x^j))_+ + \lambda \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x^i, x^j)$$

Support vector machines have sparse α 'sdecision boundary: $d(\underline{x}) = \sum_{i=1}^n \alpha_i K(\underline{x}, x^i)$

Boundary (hinge loss) depends only on the support vectors

Ex: Gaussian kernel: $K(x^i, x^j) = \exp\left\{-\frac{\|\underline{x}^i - \underline{x}^j\|^2}{2\sigma^2}\right\}$

K(x^i, x^j) measures similarity/alignment

Ex: Gaussian kernel: $K(x^i, x^j) = \exp\left\{-\frac{\|\underline{x}^i - \underline{x}^j\|^2}{2\sigma^2}\right\}$ Solve for α using gradient descent

by definition, the

$$\text{Kernel } K(\underline{x}, \underline{y}) = \phi^T(\underline{x}) \phi(\underline{y})$$

measures the similarity between 2 points

2 training data points

Process:

$$1) \text{ Initialize } w_{ij}, v_{ij} \rightarrow$$

$$2) \text{ randomly choose training sample } i$$

$$3) \text{ calculate } h_{ij}, d_i, g_i$$

$$4) \text{ Gradient descent update } V_{ij}, \text{ then } W_{ij} \text{ (deep to shallow)}$$

chain rule is key

b.3 Intro to neural networks

relate neural networks to linear classifiers

- define structure of multilayer neural network

- overview or procedure for training neural networks

- linear max margin classifier
- If we have data points that are linearly separable we can derive a classifier that splits the data based on the max difference between data points. The points that are used to find max differences are support vectors.

Support vectors

$$w = [w^T, b]^T \text{ depends only on the support vectors}$$

$$x^T w + b = 0 \quad x^T w = 0$$

$$x^T w = 0 \quad \text{distance from origin}$$

$$d(\underline{x}) = \sum_{i=1}^n \alpha_i K(\underline{x}, x^i)$$

<