

Gaussian Distribution:

$$\text{normal: } p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Unsupervised: Clustering, partitioning, k-means
K-means ++, K-medians, Vector Quantization, PCA, Parametric Density, Non-Parametric Density.

multi: $p(x) = \frac{1}{|\Sigma|^{1/2} \pi^n} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$

The covariance matrix Σ is a matrix such that each entry is $\Sigma_{ij} = \text{cov}(x_i, x_j)$. Diagonal is $\Sigma_{ii} = \text{var}(x_i)$. If Σ is diagonal, all x_i are independent.

SVM-like classifiers with a boundary: a hyperplane (a line for 2D data) that separates two classes. Support vectors are the points closest to the boundaries. γ is the margin, the distance between the boundaries and the support vectors. The parameter θ is a vector. $\theta \cdot x$ gives predictions. About θ .

- Direction of θ defines the boundary. We can choose this.
- All must be γ , as restricted by $\forall i: y_i \theta \cdot x_i \geq 1$. We cannot explicitly choose this; it depends on boundary. The restriction is turned into a cost in soft-margin SVM.

Perception picks misclassified points and updates θ just enough to classify it right. $\theta + \epsilon \cdot \theta - \theta \cdot \text{J}(x)$; outliers w/outliers show the boundary. Converges iff separable. Batch Eq: $\theta \cdot x = \sum_i y_i x_i$; $x_i = \#$ of point i was misclassified.

Hard Margin SVM maximizes the margin around the boundary. Technically, it minimizes the distance between boundary and the vectors closest to it. Batch Eqn: $\theta = \sum_i y_i x_i$; where $\|x_i\|$ such that $y_i \theta \cdot x_i \geq 1$. Converges iff separable. $x_i = 1$ if support vector.

Sometimes removing a few outliers lets us find a much higher margin, or margin at all. **Soft Margin SVM** is like the hard margin SVM but penalizes misclassifications: $\min_{\theta} \|\theta\|^2 + C \sum_i (1 - y_i \theta \cdot x_i)$; Hyperparameter C is the hardness of the margin. Low C means more misclassification but larger soft margin.

Small C	Big C
desire maximize margin $\ w\ $	Keeps most slack variables zero or small
outliers less sensitive	very sensitive
boundary more "flat"	more curvy
danger underfitting (misclassifies much training data)	overfitting (good training, bad testing)

KNN: Given an item x , find the k training items closest to x and return the result of a vote. Hyperparameter k , the number of neighbors. "Closest" can be defined by some norm (L2 by default). Overfits when k is super small.

Description Tree: Recursively split on features that yield the best split. Each tree has many nodes, which either split on a feature or a threshold or all the data the same way. Hyperparameters typically restrict complexity (max tree depth, min points at node) to penalize it. One particular one of interest d , the max # of nodes.

Overfits when tree is deep or when we are allowed to split on a very small # of items. Bagging: Make multiple trees, each with a random subset of training items. To predict, take votes from trees. Hyperparameters: # of trees, proportion of items to subset.

Random forest is bagging, except for each node, consider a random subset of features to split on. Hyperparameters: proportion of features to consider.

Support Vector Machines

$$\text{Hyperplane: } w^T x + b = 0$$

$$\text{Optimization: } \min_w \frac{1}{2} \|w\|^2 \text{ s.t. } y_i^T (w^T x + b) \geq 1, i=1, \dots, m$$

$$\text{Primal: } L_p(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i (y_i^T (w^T x + b) - 1)$$

$$\frac{\partial L_p}{\partial w} = w - \sum_i \alpha_i y^i = 0 \Rightarrow w = \sum_i \alpha_i y^i$$

$$\frac{\partial L_p}{\partial b} = -\sum_i \alpha_i y^i = 0 \quad \text{Note: } \alpha_i \neq 0 \text{ only for support vectors}$$

→ substitute derivatives into primal to get the dual

Dual:

$$L_d(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j$$

KKT says $\nabla_{\alpha} L_d(\alpha) = 0$ where $\alpha_i > 0$

Non-separable case, we add boundaries & penalize $\alpha_i \in [0, 1] + C \sum_i \alpha_i$; s.t. $y_i^T (w^T x + b) \geq 1 - \xi_i$

→ each ξ_i has upperbound of $C \Rightarrow 0 \leq \xi_i \leq C$

REGRESSION

In general, loss function consist of two parts, the loss term and regularization term

$$J(w) = \sum_i \text{loss}_i + R(w)$$

L2 regularization results in RIDGE REGRESSION used when w contains a null space

$$\min_w \|Ax - y\|^2 + \lambda \|w\|^2 \rightarrow w^* = (\lambda I + X^T X)^{-1} X^T y$$

→ works well for LARGE parameters

L1 regularization results in lasso regression used when w has a Laplace Prior

→ gives sparse results → not good for feature reduction

Logistic Regression

Classify $y \in \{0, 1\}$ to model $P(y=1|x)$

$$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\frac{\partial h_0}{\partial \theta} = \left(\frac{1}{1 + e^{-\theta^T x}} \right)^2 e^{-\theta^T x} = h_0(1 - h_0)$$

$$P(y|x; \theta) = (h_0(x))^y (1 - h_0(x))^{1-y} \rightarrow$$

$$L(\theta) = \prod_i (h_0(x_i))^y_i (1 - h_0(x_i))^{1-y_i} \rightarrow$$

$$J(\theta) = \sum_i y_i \log(h_0(x_i)) + (1-y_i) \log(1 - h_0(x_i)) \rightarrow$$

$$\nabla_{\theta} J(\theta) = \sum_i (y_i - h_0(x_i)) x_i^T = X^T (y - h_0(x)) \text{ (want max } J(\theta) \text{)}$$

Stochastic: $\theta_{+1} = \theta_0 + \eta (y_i - h_0(x_i)) x_i^T$

Batch: $\theta_{+1} = \theta_0 + \eta X^T (y - h_0(x))$

LDA and GDA

Classify $y \in \{0, 1\}$, model $p(y) = \phi^T \phi^{-1} y$ and $J(\theta_0, \theta_1, \Sigma, \Sigma) = \log(\prod_i p(x_i | y_i, \theta_i, \Sigma, \Sigma))$

$$\rightarrow \text{MLE} = \frac{1}{m} \sum_i \log(y_i^T \theta_i + 1) \text{ and } \hat{\mu}_{1, \text{true}} = \text{avg}(x_i)$$

$$\Sigma_{\text{true}} = \frac{1}{m} \sum_i (x_i - \hat{\mu}_1)(x_i - \hat{\mu}_1)^T$$

If $p(x|y)$ is multivariate gaussian w/ shared Σ , then $p(y|x)$ is logistic function. The converse is NOT true.

LDA makes stronger assumptions about the data than does linear regression.

$$h(x) = \text{argmax}_{y \in \{0, 1\}} -\frac{1}{2} (x - \mu_y)^T \Sigma^{-1} (x - \mu_y) + \log(\pi_y)$$

For GDA, it is the same except for each class has a unique covariance matrix.

$$h(x) = \text{argmax}_{y \in \{0, 1\}} -\frac{1}{2} \log(\pi_y) - \frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) + \log(\pi_y)$$

Underfitting = high bias; overfitting = high variance

Probability + Matrix Review

Bayesian Decision Theory

$$\text{Bayes Rule: } P(w|x) = \frac{P(x|w)P(w)}{P(x)} ; P(x) = \sum_i P(x|w_i)P(w_i)$$

$$P(x, w) = P(x|w)P(w) = P(w|x)P(x)$$

$$P(\text{error}) = \int_{\Omega_w} P(\text{error}|x)P(x) dx$$

$$P(\text{error}|x) = \begin{cases} P(w_1|x) & \text{if we decide } w_1 \\ P(w_2|x) & \text{if we decide } w_2 \end{cases}$$

$$0-1 \text{ Loss: } \lambda(w_i|x) = \begin{cases} 1 & \text{if } x \in \Omega_i \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Expected Loss (Risk): } R(w|x) = \sum_i \lambda(w_i|x)P(w_i|x)$$

$$0-1 \text{ Risk: } R(w|x) = \sum_i P(w_i|x) = 1 - P(w|x)$$

Generative vs Discriminative Model

Generative: Model class conditional density $p(x|y)$ for $p(y|x) \propto p(x|y)p(y)$ or model joint $p(x, y) = p(x|y)p(y)$ and marginalize to find $p(y|x) = \int_x p(x, y) dx$

Discriminative: Model condition $p(y|x)$ class conditional $p(x|y)$; posterior $p(y|x)$ prior $p(y)$; evidence $P(x)$

Loss Functions

$$\text{Binomial deviance} = \log[1 + e^{-f(x)}]$$

$$\text{minimizing function: } f(x) = \log \frac{p(y=1|x)}{p(y=0|x)}$$

$$\text{SVM hinge loss} = [1 - yf(x)]_+$$

$$\text{minimizing function: } f(x) = \text{sign}([1 - yf(x)]_+ - \frac{1}{2})$$

$$\text{Squared Error} = [y - f(x)]^2 = [1 - yf(x)]^2$$

$$\text{minimizing function: } f(x) = 2P[y=1|x] - 1$$

$$\text{"Huberized" Squared hinge loss}$$

$$= \begin{cases} 1 - yf(x)^2 & \text{if } yf(x) < 1 \\ [1 - yf(x)]^2 & \text{otherwise} \end{cases}$$

$$\text{minimizing function: } f(x) = 2P[y=1|x] - 1$$

Gradient Descent

$$\theta_{+1} = \theta_t - \eta \nabla_f(\theta_t), \text{ for minimizing}$$

$$\text{Gradients: } \frac{\partial L(w)}{\partial w} = \frac{\partial L(w)}{\partial w} = A^T \frac{\partial L(w)}{\partial w} = A$$

$$\frac{\partial L}{\partial w} = \begin{matrix} \frac{\partial L}{\partial w_1} & \dots & \frac{\partial L}{\partial w_n} \\ \frac{\partial L}{\partial x_1} & \dots & \frac{\partial L}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial w_m} & \dots & \frac{\partial L}{\partial x_m} \end{matrix} = \frac{\partial L(x)}{\partial x} = \frac{\partial L(w)}{\partial w} = B^T$$

$$\frac{\partial L}{\partial w} = (A + B^T)x = (A + B^T)x$$

Neural Networks

- Neural Networks explore what you can do by combining perceptrons, each of which is a simple linear classifier. We use a soft threshold for each activation function θ because it is twice differentiable.



$$\text{Activation functions: } \begin{aligned} \theta(s) &= \tanh(s) = \frac{e^{-s} - e^s}{e^{-s} + e^s} \\ &\Rightarrow \theta'(s) = 1 - \theta^2(s) \\ \theta(s) &= \sigma(s) = \frac{1 - e^{-s}}{1 + e^{-s}} \\ &\Rightarrow \theta'(s) = \sigma(s)(1 - \sigma(s)) \end{aligned}$$

Error Functions:

$$\text{Cross Entropy Loss: } \sum_i y_i \log(h_\theta(x_i)) + (1-y_i) \log(1-h_\theta(x_i))$$

$$\text{Mean Squared Error: } \sum_i (y_i - h_\theta(x_i))^2$$

Notations:

1. w_{ij} is the weight from neuron i in layer $(l-1)$ to neuron j in layer l . There are a_l^2 nodes in the l^{th} layer.

2. L layers, where L is the output layer and data is the 0^{th} layer.

3. $x_j = \theta(s_j)$ is the output of a neuron. Its the activation function applied to the input signal. $s_j = \sum_i w_{ij} x_i + b_j$

4. $e(x)$ is the error as a function of the weights.

The goal is to learn the weights w_{ij} . We use gradient descent, but the error is non-convex so we tend to local minima. The naive version takes $O(n^2)$.

Backpropagation: an algorithm for efficient computation of the gradient, takes $O(L)$.

$$\nabla_{\theta} e(w) = \frac{\partial e(w)}{\partial w} = \frac{\partial e}{\partial w_{ij}} = \frac{\partial e}{\partial s_j} = \frac{\partial e}{\partial x_j} = \frac{\partial e}{\partial s_j} \frac{\partial s_j}{\partial x_j} = \frac{\partial e}{\partial s_j} \frac{\partial s_j}{\partial w_{ij}}$$

$$\text{Final Layer: } \theta_L^L = \frac{\partial e(w)}{\partial w} = \frac{\partial e}{\partial w_{ij}} = \frac{\partial e}{\partial s_j} = e'(x_j) \theta_{j+1}^L(s_j)$$

$$\text{General: } \theta_i^L = \frac{\partial e(w)}{\partial w} = \sum_{j=1}^L \frac{\partial e}{\partial w_{ij}} = \frac{\partial e}{\partial s_j} \frac{\partial s_j}{\partial w_{ij}} = \frac{\partial e}{\partial s_j} \frac{\partial s_j}{\partial w_{ij}} = \sum_{j=1}^L \theta_j^L \cdot w_{ij} \cdot \theta'(s_j)$$

Process:

1. Initialize all weights w_{ij} at random

2. for $t = 1, 2, \dots, d$

• pick $n \in \{1, 2, \dots, N\}$

• Forward: compute all x_i^T

• Backwards: compute all θ_j^L

3. Update the weights $w_{ij} \leftarrow w_{ij} - \eta \cdot x_i^T \cdot \delta_j^L$

4. Iterate to the next step until you reach goal, Δ

5. Return the final weights w_{ij}

LDA vs Logistic Regression

LDA Advantages:

• LDA stable for well separated classes. Logistic regression is surprisingly unstable

• Easy and elegant for more than 2 classes.

• LDA slightly more accurate when class nearly normal, especially small n

Logistic Regression:

• More emphasis on decision boundary; always separates linearly separable points thru weighting of points.

• More robust than non-Gaussian distributions.

• Naturally fits labels between 0 and 1

Clustering (cont.)

Partitioning:

Partition the data into k mutually exclusive exhaustive groups (ie: encode $h: C \rightarrow \{1, \dots, k\}$). Iteratively reallocate to minimize some loss function. Finding the correct partitions is hard. Use a greedy algorithm called k-means (coordinate descent). Loss function is non-convex and thus we find a local minimum.

k-means: Choose clusters at random, calculate centroid of each cluster, reallocate objects to nearest centroid, then repeat. Works with spherical, well-separated clusters of similar count and volumes.

k-Means++: Initialize clusters one by one. $D(x) = \text{distance of point } x$ to nearest cluster. $P(x)$ is new cluster center $\propto D(x)^2$.

k-medians: Works with arbitrary distance/dissimilarity metric, the centers μ_k are represented by data points. Is more restrictive and thus has higher loss.

General Loss: $\sum_{i=1}^n \sum_{k=1}^K \delta(x_i, \mu_k) P(x_i) = 1 - P(x_i)$

Parametric Density Estimation

Mixture Models: Assume PDF is made up of multiple gaussians with different centers.

$$P(x) = \sum_i P(c_i) P(x|c_i)$$

objective function as log likelihood of data. Use EM to estimate this model

$$E\text{-step: } P(w_i|x) = \frac{P(w_i) P(x|w_i)}{\sum_j P(w_j) P(x|w_j)}$$

$$M\text{-step: } P(c_i) = \frac{1}{n} \sum_i P(w_i|x_i)$$

$$\mu_i = \frac{\sum_i w_i x_i}{\sum_i w_i}$$

$$\Sigma_i = \frac{1}{n} \sum_i w_i (x_i - \mu_i)(x_i - \mu_i)^T$$

Non Parametric Density Estimation

Can use histogram or Kernel Density Estimation (KDE).

$$KDE: P(x) = \frac{1}{n} \sum_i K(x - x_i)$$

The kernel K has the following properties: symmetric, Normalized

$$\int K(x) dx = 1 \text{ and } \int (x-a)^2 K(x) dx = 0$$

The bandwidth is the width of the kernel function. Too small = jagged results, too large = smoothed out results.

Principal Component Analysis

First run singular value decomposition on pattern matrix X .

1. Subtract mean from each point

2. (sometimes) Scale by variance.

3. Compute covariance $\Sigma = X^T X$

4. Compute eigenvalues/vectors of $\Sigma = V \Sigma^{-1} V^T$ (spectral theorem)

5. get back $X = X \Sigma = (XV) \Sigma V^T = U \Sigma V^T$

S contains the eigen values of transformed features. The larger the S_{ii} , the larger the variance of that feature. We want the k largest features, so we find the indices of S and we only keep these entries in U and V .

MLE: To do gaussian discriminant analysis, we must fit gaussians to the sample points and estimate the class prior probabilities.

Definition: MLE is method of estimating the parameters of a statistical model by picking the parameters that maximize the likelihood function Z .

• one way of density estimation

LiKelyhood of a Gaussian given sample points x_1, \dots, x_n

$$\text{LiKelyhood} = \prod_i P(x_i | \mu, \Sigma)$$

Then we can use gradients to find out μ and Σ

• We don't know μ exactly, so we must use the estimator μ to solve for μ

• We use the sample mean and variance for $N(\mu, \sigma^2)$

good luck! - Jill from 4:30am

ROC CURVES

Receiver Operating Characteristics
Way to evaluate classifier after it is trained by running on the test set or validation set.

We can change the posterior probability threshold for Gaussian discriminant analysis or logistic regression to trade off false positives vs false negatives.

Specificity: Horizontal distance from curve to the right (1 - false pos rate)

Sensitivity: Vertical distance from bottom to curve (1 - false neg rate)

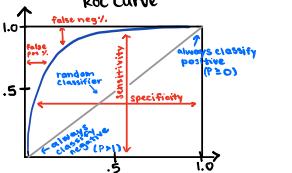
Upper right: Always classify positive

Lower left: Always classify negative

Diagonal: Random Classifiers

Rough measure of effectiveness is area under the curve. Max area is 1, for a classifier that is always right.

1/2 for random classifier.



Generative / discriminative cont.

Advantages: $A = G - P(Y|X)$ tells you probability guess was wrong, you can diagnose outliers

Disadvantages: G - often hard to estimate distributions accurately. real distributions rarely match.

Neural Networks (cont.)

Input standardization: batch norm makes sure each feature has similar statistics over a minibatch. LayerNorm makes sure each data point has similar statistics over all its features.

Skip connections: add input to result later in the network

Dropout: randomly 0 out some nodes during training and is a form of regularization encouraging network to learn diffuse rep.

Transformers:
motivated by natural language processing, since one word may not translate directly, as well as the sentences having possibly irrelevant words (ie: a, the, etc)
inductive bias - focus on pairings of words

Encoder:
full transformer includes cross-attention
similar mechanism, just querying into a different sequence

Multihead Attention:
include multiple key/query/value embedding, each is independent
Allows for learning different features

Layer Norms:
standardizes features at each position in a sentence
stabilizes training/inputs to each embedding

Transformer VIT
take each patch of the image, flatten it, and embed it with a linear layer

Now treat those embeddings like you would treat word embeddings and proceed as usual.

Can put an MLP onto some output to classify

Self Attention (transformer cont.)

each vector has dimension d_v

Attention(Q, K, V) = Softmax($\frac{QK^T}{\sqrt{d_k}}$)V

Queries - used to index information from other positions

Keys - indexed into to release the "value" from this position

Values - weighted together to form output vector

Multihead Attention:
learn different pairwise dependencies

each head is computed on its own, so we can parallelize

stack output vectors and project to desired output dimensions

↳ has dimension of nheads × d_v

project downwards w/ matrix multiplication

Bias - Variance Decomposition

Model: $X \sim D, \epsilon_i \sim D'$, $y_i = g(x_i) + \epsilon_i$; (D' has 0 mean). We fit hypothesis h (a random variable) to X , y . There are two sources of error in a hypothesis, h :

1. **Bias:** err due to inability of hypothesis h to fit g perfectly (ie: fitting a quadratic to a linear function)

2. **Variance:** err due to fitting random noise in data

Consider arbitrary point $z \in \mathbb{R}^d$, $y = g(z) + \epsilon$, $\epsilon \sim D'$

The mean comes from g and the variance comes from $\epsilon \sim E(\epsilon) = g(z)$; $\text{Var}(\epsilon) = \text{Var}(\epsilon)$

$R(h) = E[L(h(z), y)]$: We are taking a mean over the probability distribution of the hypothesis. The bias-variance decomposition of the risk function is as follows:

$$= E[(h(z) - y)^2]$$

$$= (E[h(z)] - g(z))^2 + \text{Var}(h(z)) + \text{Var}(\epsilon)$$

Visually the bias is the expected hypothesis.

Variance is the expected square distance between a random test point and the actual curve.

Consequences:

Training error reflects bias but not variance; test error reflects both

We can't precisely measure bias or variance because we don't know g and our noise model might be wrong. However, we can test learning algorithms by choosing g and making synthetic data.

Bias:

underfitting means too much bias
if h can fit g exactly, for many distributions bias $\rightarrow 0$ as $n \rightarrow \infty$

if h cannot fit g well, bias is large at most points

Adding a good feature reduces bias, adding a bad feature rarely increases it.

Bias can be 0 when hypothesis function can fit the real one.

Variance:

most overfitting caused by too much variance

For most distributions, $\text{Var} \rightarrow 0$ as $n \rightarrow \infty$

Adding a feature usually increases variance (we don't add feature unless it reduces bias more)

Variance portion of RSS (overfitting) decreases with $1/n$ sample points, increases with # features,

Irreducible Noise

Can't reduce irreducible error

Noise in test set only affects $\text{Var}(\epsilon)$. Noise in the training set only affects bias and $\text{Var}(h)$

Gaussian Discriminant Analysis cont.

Normal PDF: $f(x) = n(g(x))$, $n(x) = \frac{1}{\sqrt{(2\pi)^d}} e^{-\frac{|x-\mu|^2}{2}}$, $g(x) = (x-\mu)^T \Sigma^{-1} (x-\mu)$

$\Sigma \rightarrow R$, exponential $\Sigma \rightarrow R$, quadratic

$\Sigma^{-1} = \Sigma^{-1}$ is the covariance matrix, where eigenvalues of Σ are the variances along the eigenvectors

$\Sigma^{-1} = \sigma_i^2$

$\Sigma^{-1/2} = \Sigma^{-1/2}$ maps spheres to ellipsoids, where eigen values of Σ are the Gaussian widths/ellipsoid radii / standard deviations

$\Sigma^{-1/2} = \Sigma^{-1/2}$ is the precision matrix.

Spherically or Gaussian

$X \sim N(\mu, \Sigma)$: we want $\epsilon = f(x) \approx Z \sim N(0, I)$

$p(x) = \exp(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu))$

$p(z) = \exp(-\frac{1}{2} z^T z)$

$g(x) = (x-\mu)^T \Sigma^{-1} (x-\mu)$

Then we want $g(f(x)) = f(x)^T f(x)$

for some transformation.

$$q(x) = (x-\mu)^T \Sigma^{-1} (x-\mu)$$

$$= (x-\mu)^T V \Lambda^{-1} V^T (x-\mu)$$

$$= [(x-\mu)^T V] \Lambda^{-1} [V^T (x-\mu)]$$

$$= f(x)^T f(x)$$

↳ where $f(x) = \Lambda^{-1/2} V^T (x-\mu)$

$$\text{So } \therefore z = \Lambda^{-1/2} V^T (x-\mu)$$

divide each eigen vector direction by standard deviation in that direction

rotate the eigen basis around center

scale along axis

Another way to write this:

$$z = A^{-1}(x-\mu) \text{ for } \Sigma = AA^T$$

$$A = V \Lambda^{1/2}$$

$$\therefore A^{-1} = \Lambda^{-1/2} V^T$$

3. rotate eigen basis

4. scale along axis

Convolution Neural Network

Inductive Bias: exploiting human knowledge of a problem in order to simplify our models

Fully connected networks (MLP) can in theory learn any function, however, they have many parameters (lots of memory + run time constraints) and they require ridiculous amounts of data to learn high dimensions

Some problems have structural invariance and equivariance

Want # of parameters to be constant wrt image size

Pooling layer - used to reduce spacial dimensions

Sliding window method similar to convolution but has no parameters.

Avg pooling (take average values)

- extracts features smoothly

Max pooling (take maximum value)

- extracts most important features

of parameters = $(s+1) \times k$

k = # of filters; s = dimensions of filter (ie: $3 \times 3 \times 2 = 30 = s$)

shape of output tensor w/ k x k filter

$[H, W, K] = 1 \times ([H, W] + 2 \times pad - k) / stride$

Fully connected layers typically go on the end of models w/ reduced dimension features

Math Appendix

Vectors

Dot Product: $X \cdot Y = X_1 Y_1 + X_2 Y_2 + \dots + X_n Y_n$

Euclidean norm: $\|X\|_2 = \sqrt{X_1^2 + X_2^2 + \dots + X_n^2}$

'Normalizing' a vector means replacing X with $\|X\|_2$

$\cos(\theta) = \frac{X \cdot Y}{\|X\|_2 \|Y\|_2}$

Continuous Distributions

$P(X \in [x_1, x_2]) = \int_{x_1}^{x_2} f(x) dx$

$E[X] = \int_{-\infty}^{\infty} x f(x) dx$

$E[X^2] = \int_{-\infty}^{\infty} x^2 f(x) dx$

$\sigma^2 = E[(X - \mu)^2] = E[X^2] - \mu^2$

Linear Algebra

Basics

The column space, also called the range or span of X is

$\text{Range}(X) = \{y | y = Xv\}$

The row space is $\text{Row}(X) = \{y | y = Xv\}$

The null space, or the Kernel of X is defined as $N(X) = \{v | Xv = 0\}$

The orthogonal complement of a subspace U , is a subspace U^\perp such that $u \in U, u' \in U^\perp \Rightarrow u \cdot u' = 0$

$\text{Row}(X) = \text{Range}(X^T)$

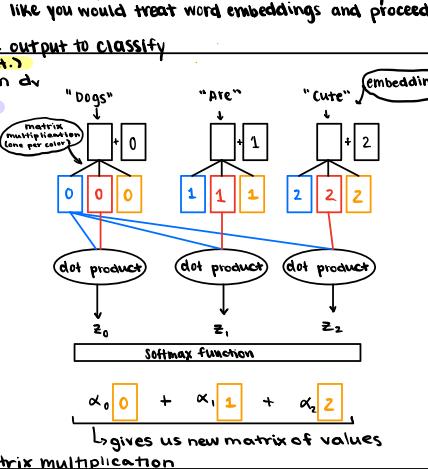
$N(X)^{\perp} = \text{Row}(X)$

$N(X^T X) = N(X)$

$\text{Row}(X^T X) = \text{Range}(X^T X)$

Rank Nullity Theorem: $\text{Rank}(X) + \text{Rank}(N(X)) = m$

$\text{Rank}(A) + \text{Rank}(B) - n \leq \text{Rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$



↳ gives us new matrix of values