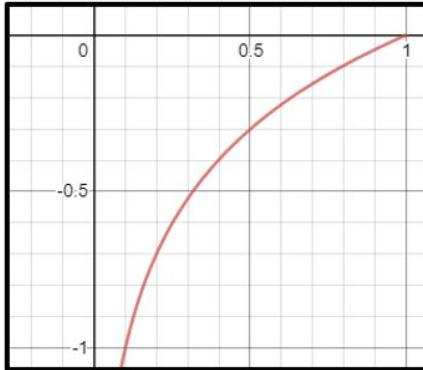
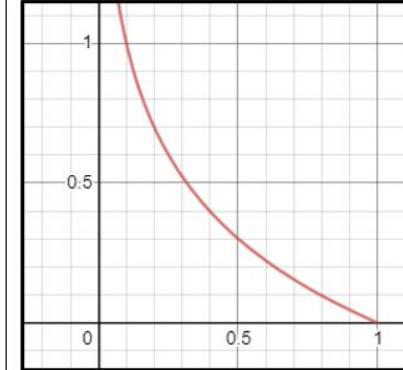
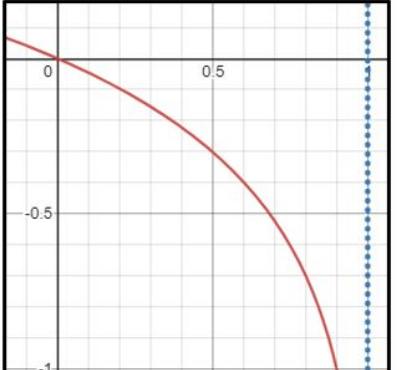
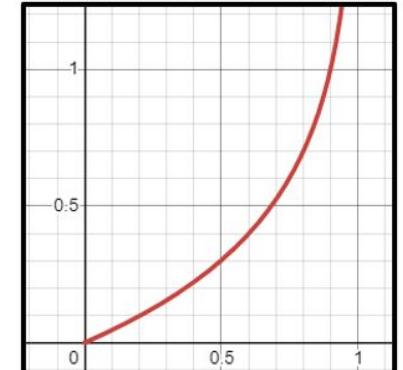


Loss Functions & Metrics

Describe and draw the properties of $\log(x)$, $\log(1-x)$, $-\log(x)$, $-\log(1-x)$

$f = \log(x)$	$f = -\log(x)$	$f = \log(1-x)$	$f = -\log(1-x)$
			
<p>f is small if x is small – they “move in the same direction”</p>	<p>f is small if x is large – they “move in opposite directions”</p>	<p>f is small if x is large – they “move in opposite directions”</p>	<p>f is small if x is small – they “move in the same direction”</p>

Complete the following confusion table & associated metrics. Are they used for classification or prediction?

		Predicted Condition			
Actual Condition	Actual Positive (P)	Predicted Positive (PP)	Predicted Negative (PN)		
	Actual Negative (N)				
Total Population					
	Matthews correlation coefficient (MCC) $= \sqrt{TPR \cdot TNR \cdot PPV \cdot NPV - \sqrt{FNRR \cdot FPR \cdot FOR \cdot FDR}}$ (Considered most balanced, single metric to evaluate confusion matrix)			Balanced Accuracy	Macro-Avg-Precision = Macro-Avg-Recall = Macro-Avg-F1 =
				Accuracy	
				F_1 Score	

- **Prediction** is free of decisions; it simply outputs a vector of probabilities. It is best used for nuanced, low signal-to-noise ratio scenarios (eg medical diagnosis, weather forecasting).
- **Classification** applies a threshold on these probabilities depending on the particular use case/costs associated -- a decision & judgement call is made. It is best used for non-deterministic, high signal-to-noise ratio scenarios where outcomes are distinct (eg character recognition)
- These metrics help with **interpreting the results of classification** algorithms, and understanding their success/failure magnitude.
- They can be defined for several classification thresholds.

		Predicted Condition			
Actual Condition	Actual Positive (P)	Predicted Positive (PP)	Predicted Negative (PN)		
	Actual Negative (N)				
True Positive (TP) "Hit" "Correctly rejected null hypothesis"	False Negative (FN) "Type II error", "Miss" "Incorrectly kept null hypothesis"	True Positive Rate (TPR) "Recall", "Sensitivity", "Power", "Completeness" $= \frac{TP}{P} = 1 - FNR$ $= P(\hat{Y} = 1 Y = 1)$	False Negative Rate (FNR) "Miss rate" $= \frac{FN}{P} = 1 - TPR$ $= P(\hat{Y} = 0 Y = 1)$		
False Positive (FP) "Type I error", "False alarm" "Incorrectly rejected null hypothesis"	True Negative (TN) "Correct rejection" "Correctly kept null hypothesis"	False Positive Rate (FPR) "Probability of false alarm" $= \frac{FP}{N} = 1 - TNR$ $= P(\hat{Y} = 1 Y = 0)$	True Negative Rate (TNR) "Specificity" $= \frac{TN}{N} = 1 - FPR$ $= P(\hat{Y} = 0 Y = 0)$		
Positive Predictive Value (PPV) "Precision", "Exactness" $= \frac{TP}{P} = 1 - FDR$ $= P(Y = 1 \hat{Y} = 1)$	False Omission Rate (FOR) $= \frac{FN}{P} = 1 - NPV$ $= P(Y = 1 \hat{Y} = 0)$				
False Discovery Rate (FDR) $= \frac{FP}{P} = 1 - PPV$ $= P(Y = 0 \hat{Y} = 1)$	Negative Predictive Value (NPV) $= \frac{TN}{N} = 1 - FOR$ $= P(Y = 0 \hat{Y} = 0)$				
Total Population $= P + N = (TP + FN) + (TN + FP)$ $= PP + PN = (TP + FP) + (TN + FN)$	Prevalence "Proportion affected" $= \frac{P}{P + N}$	Accuracy $= \frac{TP + TN}{P + N}$	F_1 Score $= \text{harmonic mean}(precision, recall)$ $= \frac{2PPV * TPR}{PPV + TPR}$		
	Matthews correlation coefficient (MCC) $= \sqrt{TPR \cdot TNR \cdot PPV \cdot NPV - \sqrt{FNRR \cdot FPR \cdot FOR \cdot FDR}}$ (Considered most balanced, single metric to evaluate confusion matrix)	Balanced Accuracy $= \frac{TPR + TNR}{2}$		(Avg metrics of current and flipped pos/neg labels; allows looking at P/R at ALL classes)	
				Macro-Avg-Precision = $(PPV + NPV)/2$	
				Macro-Avg-Recall = $(TPR + TNR)/2$	
				Macro-Avg-F1 = $(\frac{2PPV * TPR}{PPV + TPR} + \frac{2NPV * TNR}{NPV + TNR})/2$	

Intuitively, what is precision and recall? When should you pick one over the other?

- **Positive Predictive Value = Precision** = $PPV = \frac{TP}{(TP + FP)}$
 - Is about **exactness**: measures ability to reliably reject non-relevant documents
 - Punishes False Positives
 - Important if costly for a non-relevant item to get mistakenly accepted (often occurs if “there are plenty of fish in the ocean”, such as for political elections or job applications)
- **True Positive Rate = Recall** = $TPR = \frac{TP}{(TP + FN)}$
 - Is about **completeness**: measures ability to reliably find all relevant documents
 - Punishes False Negatives
 - Important if costly for relevant item to get missed (such as in safety-critical applications like autonomous driving or medicine)

What is the harmonic mean, and why is it used in the F1 score?

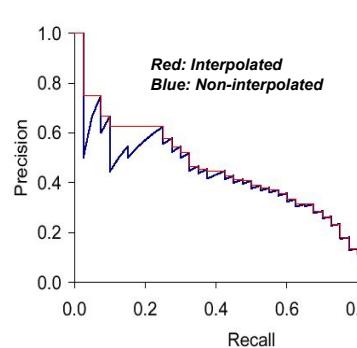
- In general, for averages to be valid, you need the denominator values to be in the same scaled units
- This exactly what the harmonic mean does: it is the reciprocal of the arithmetic mean of the reciprocals of the given set of observations
 - Precision and recall both have true positives in the numerator, and different denominators. To average them it really only makes sense to average their reciprocals, thus the harmonic mean.

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{\text{tp}}{\text{tp} + \frac{1}{2}(\text{fp} + \text{fn})}$$

$$\left(\frac{x_1^{-1} + \dots + x_n^{-1}}{n} \right)^{-1} \quad \left(\frac{1^{-1} + 4^{-1} + 4^{-1}}{3} \right)^{-1} = \frac{3}{\frac{1}{1} + \frac{1}{4} + \frac{1}{4}} = \frac{3}{1.5} = 2$$

How do you compute precision/recall @ K? What is a precision-recall curve? What is the difference between average precision (AP) and mean average precision (mAP)?

- Suppose we have some **ranked ordering/retrieval** of the positive predictions. Then:
 - **Precision @ K** = (# relevant out of top K) / K = (TP in K) / K
 - Can go up or down as you increase K
 - **Recall @ K** = (# relevant out of top K) / (# relevant in total) = (TP in K) / TP
 - Monotonically increasing; cannot be 1 until K >= # relevant
- Note that true/false negatives are not taken into account; we just want the positives to be ranked highly
- **Precision-Recall (PR) Curves** plots the precision against the recall at all values of K.
 - In general they tend to be inversely proportional (e.g. higher recall -> lower precision)
 - Also, PR curves are more useful when the positive class is rare/interesting. They ignore True Negatives.
 - When plotting these curves, we can **interpolate** and use the best precision for a level of recall R or greater.
 - In some literatures, this is standard (the justification is that almost anyone would be prepared to look at a few more documents if it would increase the percentage of the viewed set that were relevant).
 - However, some view this as being too optimistic, and do not perform the interpolation.
- **Average Precision (AP)** is calculated by finding the area under the precision-recall curve (can be either interpolated or not)
 - For example, scikit-learn only gives the non-interpolated version (see equation, where there are n thresholds)
- **Mean Average Precision (mAP)** is the AP averaged across a dataset (eg, an AP for each class, AP for each IoU thresholds, etc)



$$AP = \sum_n (R_n - R_{n-1}) P_n$$

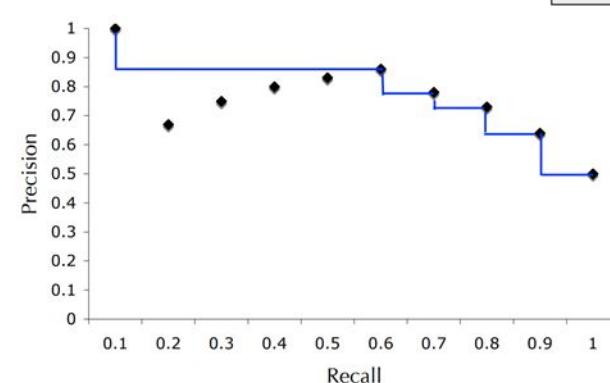
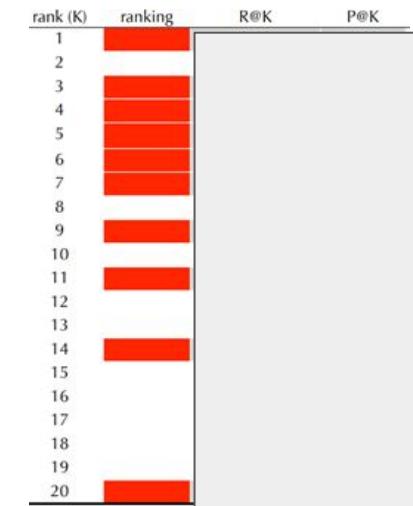
$$mAP = \frac{1}{N} \sum_{i=1}^N AP_i$$

Calculate the accuracy, precision, recall, and F1 score for the following binary confusion matrix.

Complete the following P/R@K table assuming there are 10 true positives, and draw the PR curve.

Predicted	Actual		
	Bipolar	Not Bipolar	Total
Bipolar	10	11	21
Not Bipolar	9	52	61
Total	19	63	82

- Accuracy
 - 62/82
- Precision
 - 10/21
- Recall
 - 10/19
- F1 Score
 - $(2 * (10/21) * (10/19)) / ((10/21) + (10/19))$



How is mAP calculated in the context of retrieval vs object detection?

- **Retrieval**
 - The “standard” case; check the top K retrieved images in the test set for a retrieval query (eg, an object class) and compute the precision/recall for each K by considering if they are relevant to your class. Thus, you get the PR curve, and calculate the AP.
- **Object detection**
 - For a given object class, precision/recall value calculated for each detection (after sorting by descending confidence) over every image the test dataset
 - True/false positive detections are based on a chosen IoU threshold
 - Plotting gives a class PR curve from which you can calculate the AP (by taking the AUC)
 - This is repeated for every class
- In both cases, mAP can be computed by averaging the APs for each class. For object detection, this might also be averaged over different IoU thresholds.

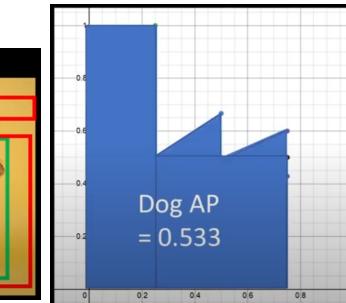
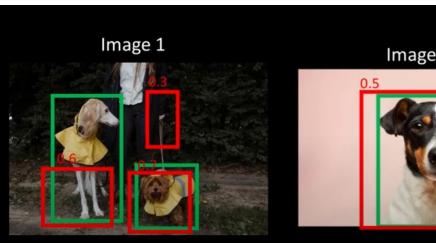
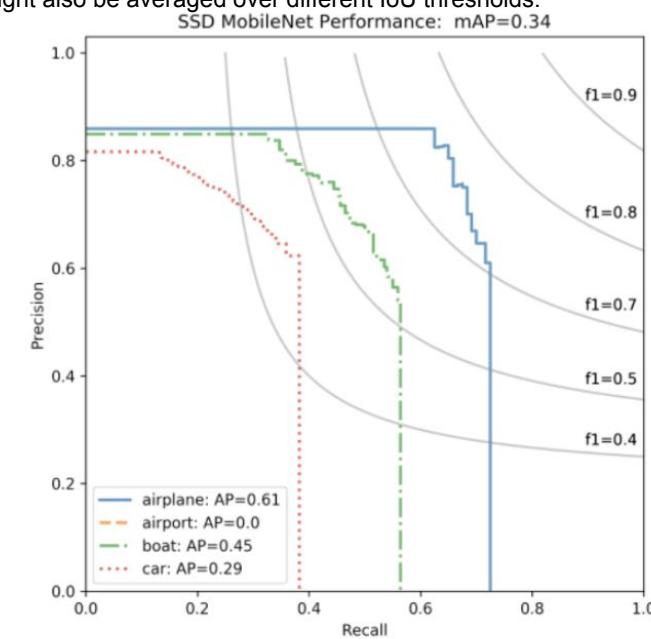


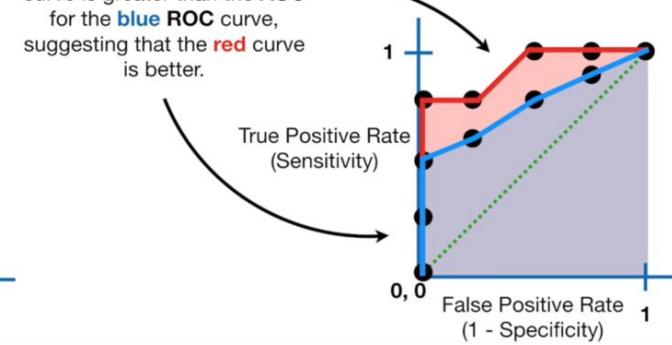
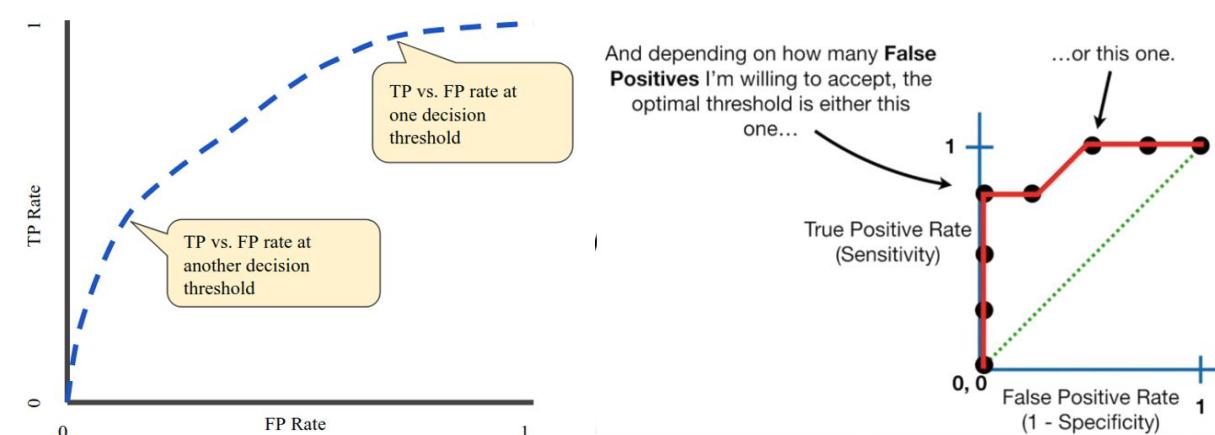
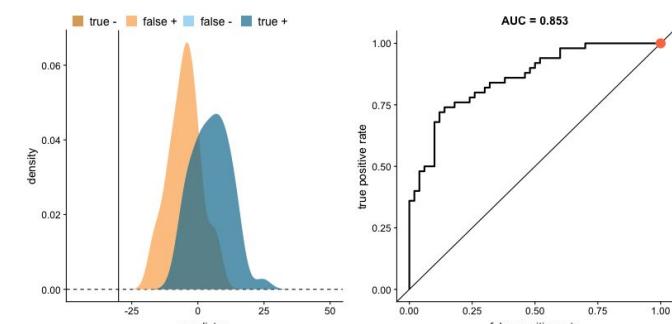
Image	Confidence	TP or FP	Precision	Recall
Image 3	0.9	TP	1 / 1	1 / 4
Image 3	0.8	FP	1 / 2	1 / 4
Image 1	0.7	TP	2 / 3	2 / 4
Image 1	0.6	FP	2 / 4	2 / 4
Image 2	0.5	TP	3 / 5	3 / 4
Image 1	0.3	FP	3 / 6	3 / 4
Image 3	0.2	FP	3 / 7	3 / 4



Precision-recall curve for SSD model for 4 object classes, where IoU threshold is 0.5. Van Etten, A. (2019, January).

What is an ROC curve, and how is it similar/different from the PR curve?

- An ROC curve is created by plotting the true positive rate (TPR/Recall) against the false positive rate (FPR) at various classification threshold settings
 - TPR/Recall:** $TP/(TP+FN) = TP / P$
 - FPR:** $FP/(TN + FP) = FP / N$
 - In general, **TPR and FPR are proportional**
 - If you want to increase the TPR, you inadvertently increase the FPR
 - If you want to decrease the FPR, you inadvertently decrease the TPR
- The AUC (area under the curve) of the ROC curve is a summary of the performance of a method
- Unlike PR curves, ROC curves take into account true negatives



How do you know which classification metrics to use?

Advice on when to use which metrics:

- Only care about correct classifications of **actual positives xor actual negatives**, especially if the class you care about is in the minority:
 - Use **precision**, **recall**, **Average precision**, and/or **PR-curves**, which are tailored for rare events
 - If you only care about the positive cases, you're all set
 - If you only care about the negative cases, simply flip the classes (treat the negative class as positive, and vice versa)
- Care about correctly classifying **both actual positives and negatives**:
 - Use **balanced accuracy**, **MCC**, **macro-averaged precision/recall/F1**, and/or **ROC curves**
 - ROC is best for a balanced scenario; can be misleading for imbalanced scenarios
 - Macro-averaged precision/recall/F1 allows you to look at precision/recall/F1 at ALL your classes, and aggregates it

Define Shannon information, entropy, cross entropy, KL/JS Divergence. What are their relationships?

Note: the following applies to discrete random variables, though there are extensions to the continuous case.

- **Shannon Information $I(X=x)$:** Intuitively quantifies the level of “surprise” of a probabilistic realization, by taking the negative log of its probability
 - It’s an inverse relationship; the higher the likelihood/probability of that realization, the lower the shannon information
 - As p approaches 0, the entropy approaches infinity; if $p = 1$, entropy = 0.
- **Shannon Entropy $H(X)$:** The average amount of uncertainty of a R.V., i.e. the expected value of the shannon information
 - Higher (lower) entropy means that the R.V. on average has more (less) uncertainty/randomness.
 - Maximized when the distribution is multinoulli uniform; minimized when an event has probability 1 while all others is 0
 - If $b=2$, we can interpret $H(X)$ as the lower bound for the average number of bits needed to encode a piece of data ($X=x$). Therefore, besides a notion of uncertainty entropy also represents a way to **quantify the “amount of information” in the data**.
 - For example, if $H(X)=0$, it’s the same value over and over again; no need to encode.
 - If $H(X)$ is moderate, there are some underlying patterns and tendencies that allow for some optimizations.
 - If $H(X)$ is high, there are many possibilities; need many bits to fully capture data.
 - Suppose we have a fair n -sided die. Then, the entropy becomes $H(X)=-\log_2(1/n)=\log_2(n)$, which intuitively is the number of bits needed to record outcomes of a dice roll.
- **Cross Entropy $H(P,Q)$:** Intuitively, average number of **total** bits required to encode data coming from a source with “true” distribution P when we use model Q .
 - If P is a fixed reference distribution, then the cross entropy reaches its minimum $H(P)$ when $Q = P$ (from Gibbs’ inequality).
 - This is how the cross entropy loss is formulated, usually where P is a one-hot ground truth.
 - It is the same as the **negative log-likelihood**; their only difference is the way they’re interpreted
 - Also known as the **log loss**
- **KL Divergence $D(P||Q)$:** average number of **extra** bits required to encode data coming from a source with “true” distribution P when we use model Q .
 - Interpretation 2: expected log likelihood ratio of P and Q which characterizes how much more likely x is drawn from P vs Q . If $P = Q$, then it is 0. Specifically, $E_p [\log(P/Q)]$.
 - The cross entropy of P and Q is equal to the KL divergence of P and Q plus the entropy of P . If P is fixed (e.g. doesn’t change with time), then minimizing the cross entropy is the same as minimizing the KL divergence up to a constant.
- **JS Divergence $JSD(P||Q)$:** Symmetric version of the JS divergence using a mixture distribution of P and Q .

$$I_X(x) := -\log_b[p_X(x)]$$

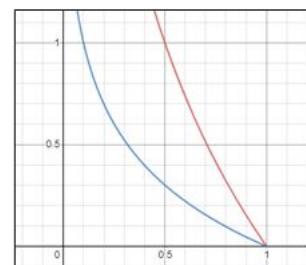
$$H(X) = E[I(X)] = - \sum_i p_X(x_i) \log_b p_X(x_i)$$

$$H(P, Q) = - \sum_{x \in \mathcal{X}} p(x) \log(q(x))$$

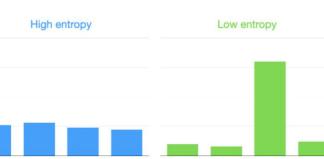
$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

$$H(P, Q) = H(P) + D_{KL}(P \parallel Q)$$

$$\begin{aligned} JSD(P \parallel Q) &= \frac{1}{2} D_{KL}(P \parallel M) + \frac{1}{2} D_{KL}(Q \parallel M) \\ M &= \frac{1}{2} (P + Q) \end{aligned}$$



Shannon information: red is $b = 2$, blue is $b = e$



$$H(P, P) = H(P)$$

$$D_{KL}(P \parallel P) = 0$$

Why do we use KL divergence for knowledge distillation rather than cross entropy, as we normally use for classification?

- Usually for classification, the target (label) distribution is one-hot, so the entropy part of the cross-entropy loss is 0
 - So in this case cross-entropy is equivalent to KL divergence; it would also be correct to say we're using the KL divergence as a loss function
 - However, the cross-entropy formulation is a bit faster to compute usually in many libraries
- For distillation, we want to compare the student's predicted label distribution to the teacher's, which is not one-hot.
 - This means the entropy portion will be non-zero; it is an additive constant
 - Since the entropy portion does not depend on the network parameters, this doesn't affect things from an optimization standpoint
 - However, it leads to the “optimal” loss being non-zero (i.e. only contributions from entropy portion) and constantly fluctuate depending on the batch, which is not as intuitive or useful.

$$H(P, Q) = - \sum_{x \in \mathcal{X}} p(x) \log(q(x))$$

$$H(X) = E[I(X)] = - \sum_i p_X(x_i) \log_b p_X(x_i)$$

Entropy (0 in the one-hot case)

$$H(P, Q) = H(P) + D_{KL}(P \parallel Q)$$

Cross Entropy

$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

KL Divergence

$$- \sum_{i=1}^N \sum_{j=1}^C y_{ic} \log (f(x_i)_c)$$

What is the MSE loss, and how does it compare to the cross entropy loss?

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$\text{MSE}(\hat{\theta}) = \mathbb{E}_{\theta} [(\hat{\theta} - \theta)^2]$$

- There are two versions of MSE: one for predictions over a labeled dataset, and one for an estimator
- In general you would use MSE for regression problems, while CE is for classification problems.
 - MSE takes into account the whole vector
 - CE only considers the entry which corresponds to the ground truth class (which is the one effectively used for classification, in one-hot cases) or is weighted by G.T. probability (in non-one-hot cases). It also makes it so that the edge values of 0 and 1 are of special importance.
- MSE is equivalent to taking the squared L2 norm between vectors

I would like to show it using an example. Assume a 6 class classification problem.

Assume, True probabilities = [1, 0, 0, 0, 0, 0]

Case 1: Predicted probabilities = [0.2, 0.16, 0.16, 0.16, 0.16, 0.16]

Case 2: Predicted probabilities = [0.4, 0.5, 0.1, 0, 0, 0]

The MSE in the Case1 and Case 2 is **0.128** and **0.1033** respectively.

Although, Case 1 is correctly predicting class 1 for the instance, the loss in Case 1 is higher than the loss in Case 2.

$$MSE(X, Y) = \frac{1}{n} \|X - Y\|_2^2 = \frac{1}{n} (\sqrt{\sum_i (X_i - Y_i)^2})^2 = \frac{1}{n} \sum_i (X_i - Y_i)^2$$

Define RMSE (Root Mean Squared Error) and MAE (Mean Absolute Error). When should you use each?

- RMSE has the benefit of penalizing large errors more so can be more appropriate in some cases
- for example, if being off by 10 is more than twice as bad as being off by 5. But if being off by 10 is just twice as bad as being off by 5, then MAE is more appropriate.

$$\text{MAE} = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

Experimental Designs and Paradigms

What is the “No Free Lunch Theorem”?

- States that "**When averaged across all possible situations, every algorithm performs equally well**"
- **For every case where an algorithm works, I could construct a situation where it fails terribly**
- Implies that **assumptions are where your power comes from**

What is the exploration-exploitation trade-off?

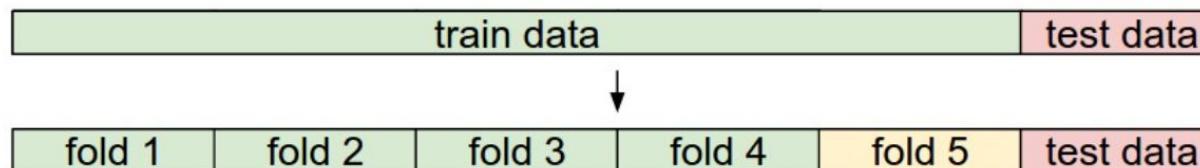
- A fundamental dilemma whenever you learn about the world by trying things out. The dilemma is between choosing:
 - Exploration
 - Something you aren't sure about and possibly learning more
 - Involves research, variation, risk taking, experimentation, discovery, and innovation
 - Exploitation
 - What you know and getting something close to what you expect
 - Involves refinement, efficiency, selection, implementation, and execution
- Has implications for business, science, engineering/research, etc
- Also used in reinforcement learning

What is k-fold cross-validation, and what are its pros and cons?

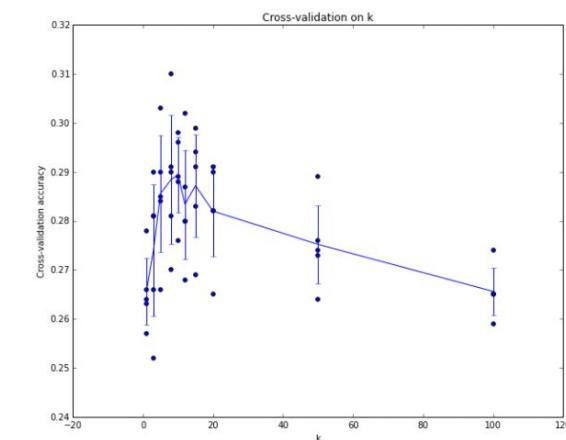
For example, in 5-fold cross-validation, we would split the training data into 5 equal folds, use 4 of them for training, and 1 for validation. We would then iterate over which fold is the validation fold, evaluate the performance, and finally average the performance across the different folds.

Pros: hyperparameter tuning becomes more reliable and less noisy

Cons: computationally expensive, linear to the number of folds

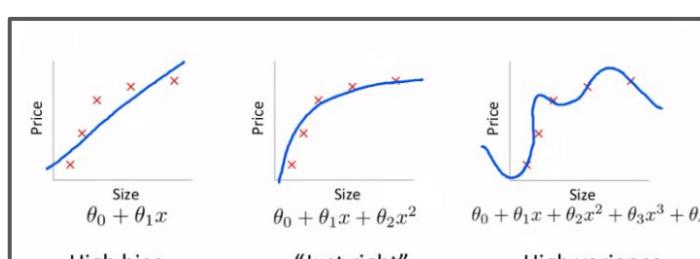


Common data splits. A training and test set is given. The training set is split into folds (for example 5 folds here). The folds 1-4 become the training set. One fold (e.g. fold 5 here in yellow) is denoted as the Validation fold and is used to tune the hyperparameters. Cross-validation goes a step further and iterates over the choice of which fold is the validation fold, separately from 1-5. This would be referred to as 5-fold cross-validation. In the very end once the model is trained and all the best hyperparameters were determined, the model is evaluated a single time on the test data (red).

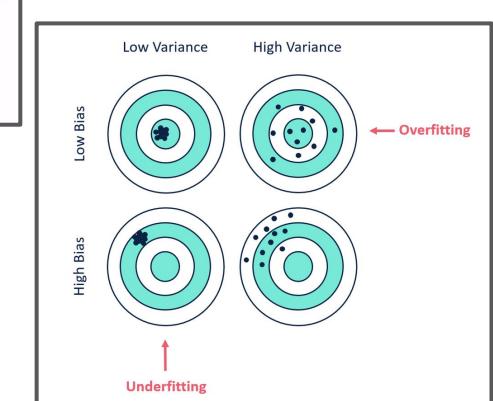
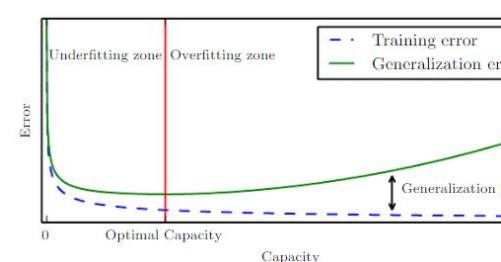
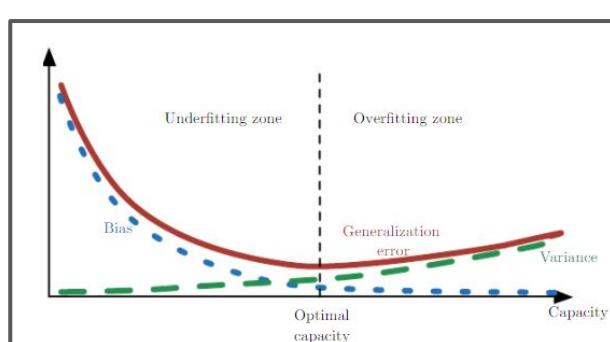


What is bias and variance? Explain their relationship, mathematical definition, and application in ML.

- **Bias** $\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta$.
 - The bias for a parameter estimator is the difference between its expected value and the true underlying parameter.
 - For ML: **on average how correct/high quality** your trained parameters will be to a true “oracle”. It should be **asymptotically unbiased**; as the number of training sample approaches infinity, the model should be unbiased. Otherwise, there’s **systematic error**.
- **Variance** $\text{Var}(\hat{\theta})$
 - Given a parameter estimator, variance characterizes how **consistent** it would be given different training datasets sampled from the underlying data-generating process.
 - For ML: if you resampled for different training datasets, would those trained parameters be similar to each other (low variance), or would it overfit to each training dataset sample, leading to significantly different parameters (high variance)?
- **Mean Squared Error**
 - The MSE for an estimator is its expected squared difference between the true underlying parameter. It can be shown that it decomposes to the bias and variance.
- **Bias variance trade-off & inverse relationship**
 - In general, high variance = overfitting and high bias = underfitting.
 - Decreasing bias tends to require model complexity leading to increased variance; decreasing variance requires regularization leading to increasing bias.



$$\begin{aligned} \text{MSE} &= \mathbb{E}[(\hat{\theta}_m - \theta)^2] \\ &= \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m) \end{aligned}$$



What is the “manifold hypothesis”?

- The Manifold Hypothesis is a heuristic that states that **real-world high-dimensional data lie on low-dimensional manifolds** embedded within the high-dimensional space.
- It explains why machine learning techniques are able to find **useful features** and produce accurate predictions from datasets that have a potentially **large number of dimensions** (variables)
- The mathematical basis for this is that a lot of the dimensions are highly **correlated**, making them redundant. In other words, the **extrinsic dimensionality of a dataset often exceeds intrinsic dimensionality of an underlying phenomena**.

Explain the nature of dataset vs feature vs model quality.

In order to be successful in machine learning, at least one of the following three must be sufficiently “good”:

- The original dataset is of high quality, that is, there is a high signal to noise ratio
- Features extracted from the dataset are discriminative and unambiguous with regards to classes. That is, they lend themselves easily to simple classification techniques and simple decision boundaries
- The model is of high quality -- assumptions taken are generally correct/reasonable, and it does not lend itself to overfitting (eg via regularization).

There are trade-offs and relationships for this; for example, if the features are extremely discriminative, then the model can be very simple and can even contain false assumptions, yet the final results are good. Meanwhile, if features extracted are not ideal, then the assumptions taken with the model become more important.

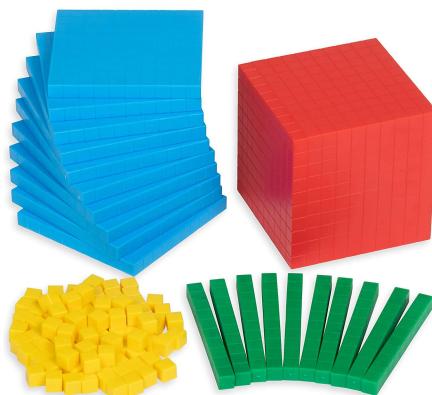
For example, a neural net learns a function such that the features are of very high quality (the class conditional distributions are approximately distributed through multivariate gaussians). Then, it can be subsequently be classified with a simple linear classifier

As a thought experiment we can take the “quality” each of these variables “to infinity”:

- Suppose the training dataset is “perfect”, and contains all the points that would be needed at test time. Then, the model’s job is simple -- just memorize the training dataset and perform lookups.
- Suppose the features are literally analogous to the class labels (perhaps in a trivially different form). Then, the model’s job is extremely trivial -- just perform a simple class-way mapping.

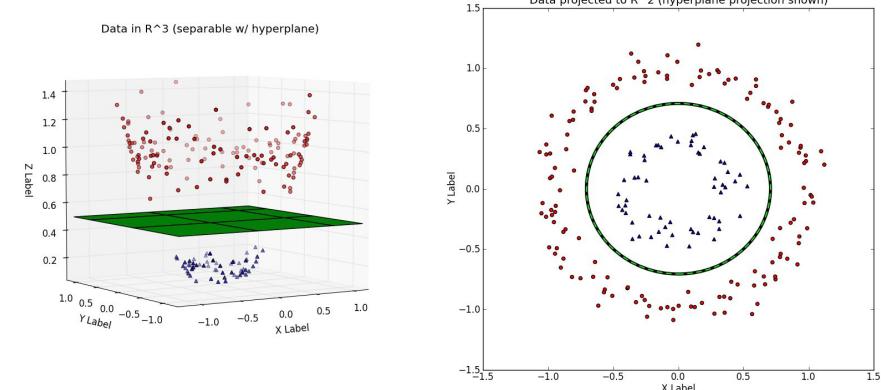
What is the “curse/blessing of dimensionality”?

- The **curse of dimensionality** refers to the problem that, when the dimensionality increases, the **volume of the space increases so fast** that the available data become **sparse**.
 - This sparsity is problematic for any method that requires **statistical significance**.
 - In order to obtain a statistically sound and reliable result, the amount of data needed to support the result often grows **exponentially** with the dimensionality.
- The blessing of dimensionality states that in some cases, there are some advantages to working with high dimensional data
 - Data becomes easier to linearly separate (eg, what kernels do)



Example of the curse of dimensionality.

- To capture 50% of 1D 10×1 grid, need 5 data points
- To capture 50% of 2D 10×10 grid, need 50 data points
- To capture 50% of 3D $10 \times 10 \times 10$ grid, need 500 data points



Example of the blessing of dimensionality; points become separable if you add a third dimension representing the radius from the middle.

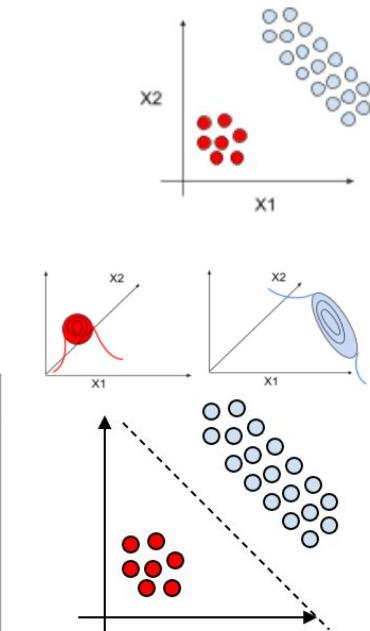
Compare/contrast generative vs discriminative methods.

The fundamental difference is that:

- **Discriminative models** learn the (hard or soft) **boundary** between classes
- **Generative models** model the **distribution** of individual classes

For example, the data to the right has two features X_1, X_2 and two classes, red and blue.

- A **generative** model would explicitly model $P(X, Y)$; in the most direct way, you would have the pdfs for $P((x_1, x_2), \text{red})$ and $P((x_1, x_2), \text{blue})$, perhaps using a gaussian mixture model
 - Classification on X then be performed by arg max over Y of $P(X, Y)$
- A **discriminative** model would implicitly or explicitly model $P(Y|X)$; given features X , you would have probability of it being red or blue.
 - We don't have $p(X)$; otherwise, we would be in the generative scenario since $P(X, Y) = P(Y|X)P(X)$
 - Classification done by arg max over Y of $P(Y|X)$



Generative	Discriminative
<ul style="list-style-type: none"> • Learns $P(X, Y) = P(X Y)P(Y)$ (implicitly or explicitly) • Models the distribution of individual classes • Less direct towards classification • More probabilistic modeling involved, more assumptions taken about the underlying distribution • Can require more assumptions • Applies to both supervised and unsupervised settings • Considers all points before modeling distribution • Access to the joint distribution; ability to draw samples from feature distribution given labels. • Offers a rich representation of the independence relationships in the dataset • Generally good for things like outlier detection • Examples: GMMs, Bayesian Networks (e.g. Naive Bayes), GANs, VAEs 	<ul style="list-style-type: none"> • Learns $P(Y X)$ (implicitly or explicitly) • Learns the decision boundary between classes • More direct way toward classification; the idea is to avoid solving an easier problem by first solving a harder problem as an intermediate step. • Fewer assumptions are taken about the data • Usually inherently supervised in nature • Generally, only considers points near decision boundary • Tends to produce higher classification accuracy • Offers a rich representation of the decision boundaries in the dataset • Examples: K-NN, Logistic regression, SVMs, Neural Networks, Random Forest

Compare/contrast classification vs regression.

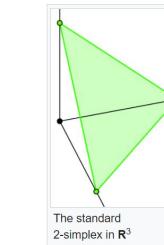
In statistical modeling and ML, we commonly want to estimate and model the relationships between a dependent variable ("labels") and one or more independent variables ("features", "inputs"). There are two settings:

- **Regression -- labels are continuous**
 - Generally, it's used for two distinct purposes:
 - In ML, it can be used to predict values given inputs
 - It can also be used to infer causal relationships between independent and dependent variables.
 - Train time evaluation is usually with MSE
 - Test time evaluation is generally with MSE, RMSE, MAE (mean absolute error)
- **Classification -- labels are discrete**
 - With the model, we can predict group categorizations. The focus is on robust decision boundaries.
 - Loss function for training: cross-entropy
 - Evaluation at test time: metrics derived from TPs, FPs, TNs, and FNs in a confusion matrix, such as accuracy (aka 0-1 loss), precision, recall, F1 scores
 - If a probability is outputted, then our output will lie on a standard probability simplex:

$$\{x \in \mathbb{R}^k : x_0 + \dots + x_{k-1} = 1, x_i \geq 0 \text{ for } i = 0, \dots, k-1\}$$

- Δ^0 is the point 1 in \mathbb{R}^1 .
- Δ^1 is the line segment joining (1, 0) and (0, 1) in \mathbb{R}^2 .
- Δ^2 is the equilateral triangle with vertices (1, 0, 0), (0, 1, 0) and (0, 0, 1) in \mathbb{R}^3 .
- Δ^3 is the regular tetrahedron with vertices (1, 0, 0, 0), (0, 1, 0, 0), (0, 0, 1, 0) and (0, 0, 0, 1) in \mathbb{R}^4 .

In many cases, regression models can be converted into classification models, and vice versa, with some small modifications.



- **Regression -> Classification**
 - Would generally involve some kind of discretization of regression model outputs, eg with some decision rule which involves some thresholding
 - softmax, sigmoid is often used to help normalize values to a range [0,1]
- **Classification -> Regression**
 - Some classification models internally have a confidence score or probability prediction for each class, before outputting the class label with maximum probability. Thus, to turn it into regression the model could just output the probabilities. We are essentially outputting soft labels instead of hard labels.
 - However this is not possible if there is no probability, e.g. SVMs are pure classification

What are the different levels of supervision that a ML model can have? Give some examples of some prevalent methods.

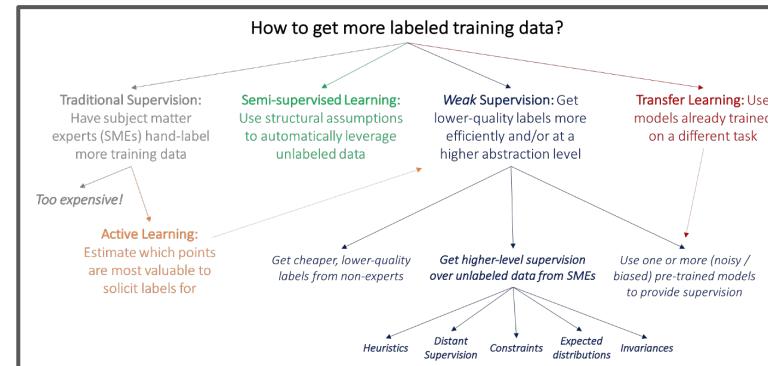
When should you consider which to use?

In machine learning, given a dataset of input features $\{X_i\}$ and labels $\{Y_i\}$. The nature of these sets lead to the following learning paradigms:

- **Fully Supervised**
 - All labels are provided for the training set. Each X_i has a corresponding Y_i , which is assumed to be ground truth.
- **Semi Supervised**
 - Some of the labels are provided for a strict subset of $\{X_i\}$. That is, only some X_i have a corresponding Y_i , which is assumed to be ground truth.
 - One technique is pseudolabeling -- predicting labels for the unlabeled features, and then training on them. This is called **bootstrapping**.
- **Weakly Supervised**
 - Noisy, limited, or imprecise sources are used to provide supervision signal for labeling training data. Each X_i has a corresponding Y_i , but Y_i could be noisy or inaccurate.
 - Some techniques are outlier detection, temporal smoothing, and modeling the noise.
- **Self Supervised**
 - Labels are generated automatically, using the inherent structure of the unlabeled training samples.
 - Examples include: Autoencoders, inpainting, superresolution. $\{Y_i\}$ is the empty set, {}.

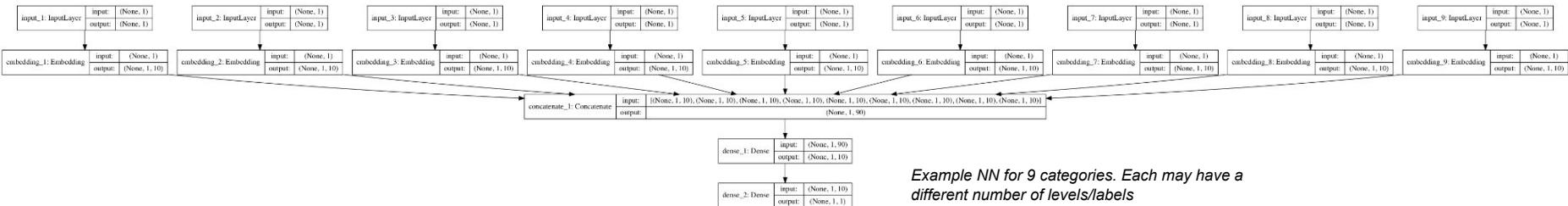
The usage generally depends on the type of data you have available and the cost of annotation.

- Semi, weakly, and self supervision can provide more training data when fully supervised from experts would otherwise be too costly.
- Another solution to this issue of obtaining labels is through transfer learning; training a model fully supervised with another dataset, and then fine tuning it fully supervised with your smaller target dataset.



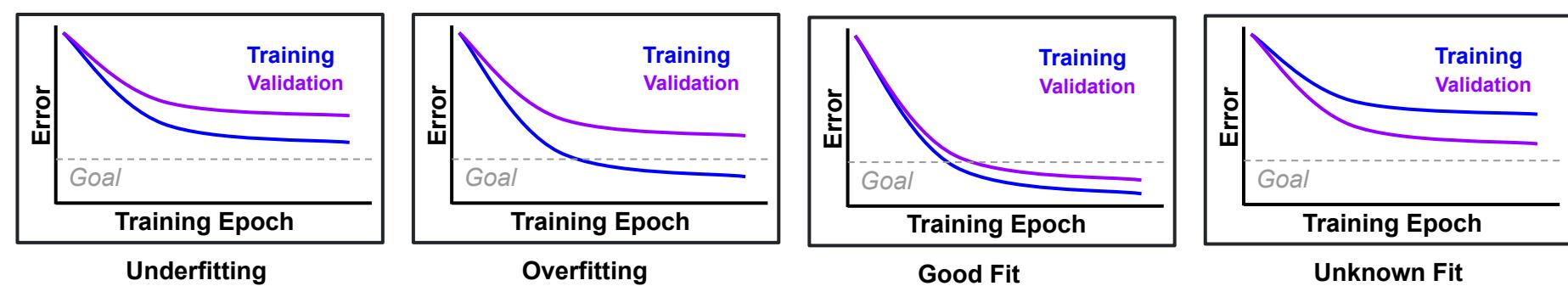
Explain & compare 3 ways to encode categorical data for NNs.

1. **One Hot Encoding:** Each label in a category mapped into a binary vector.
 - a. Simple, can always use
 - b. Sparse, and can be space inefficient, has no semantic meaning
2. **Integer Encoding:** Each label in a category mapped into an integer.
 - a. Space efficient
 - b. Only really makes sense for ordinal (ordered) data, not nominal (nonordered)
3. **Categorical Learned Embeddings:** Learn an embedding for each label in a category; use that at test time to make predictions.
 - a. Effectively tries to convert nominal data to an “ordinal-like” form. Learn embeddings can also be used for downstream tasks or post-hoc analysis
 - b. Can be more compact than one-hot encoding
 - i. For categories C1,...,Ck, need $|C_1| \times \dots \times |C_k|$ n-dimensional embeddings. This can be thought of as a matrix of learned embeddings for each category
 - c. More complicated/time consuming to train, implement, and run



Describe & illustrate the four categories of a model's performance: underfitting, overfitting, good fit, and unknown fit.

- **Underfitting** – Validation and training error high
 - Ways to mitigate: Make model bigger, reduce regularization, change model closer to SoTA, tune hyperparams, add features
- **Overfitting** – Validation error is high, training error low
 - Ways to mitigate: Collect more data, normalization, data aug, regularization, early stopping, reduce model size
- **Good fit** – Validation error low, slightly higher than the training error
- **Unknown fit** - Validation error low, training error 'high'
 - Counterintuitive to how machine learning works. The essence of ML is to predict the unknown; If you are better at predicting the unknown than what you have 'learned', then something is off
 - Could be due to bugs, or a validation set easier than the training set.
 - However, may be due to training-only regularization (e.g. more aggressive data aug or dropout)



At a high level, what is empirical risk minimization (ERM)?

- **Risk** is the theoretical expected loss over all possible x, y , for a given model
- **Empirical risk** is the expected loss over an empirical dataset of $\{x_i, y_i\}$
- **ERM** is a core principle of ML, and assuming that your dataset is sampled from the true underlying distribution, states that minimizing error on that dataset will empirically minimize the true risk (what you actually want). This comes from the law of large numbers (LLN).

$$R(h) = \mathbf{E}[L(h(x), y)] = \int L(h(x), y) dP(x, y)$$

$$R_{\text{emp}}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i)$$

What's the difference between parametric and non-parametric models? Provide some examples of both.

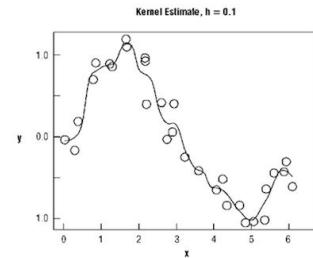
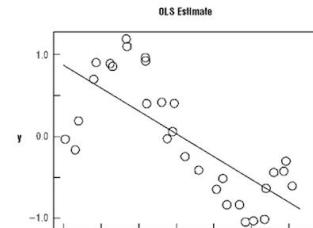
- **Parametric:**

- Number of parameters are bounded, regardless of the amount of data (sample size)
- Can be less “flexible” but less prone to overfitting
- Examples: **LDA, naive bayes, linear/logistic regression, NNs, linear SVM**

- **Non-Parametric:**

- Amount of info that parameters capture grows as the data grows
 - Effective number of params can grow with the sample size
- Can be more “flexible” but more prone to overfitting
- Examples: **K-NN, decision trees, kernelized SVM**

	Classif/regr	Gen/Discr	Param/Non
Discriminant analysis	Classif	Gen	Param
Naive Bayes classifier	Classif	Gen	Param
Tree-augmented Naive Bayes classifier	Classif	Gen	Param
Linear regression	Regr	Discrim	Param
Logistic regression	Classif	Discrim	Param
Sparse linear/ logistic regression	Both	Discrim	Param
Mixture of experts	Both	Discrim	Param
Multilayer perceptron (MLP)/ Neural network	Both	Discrim	Param
Conditional random field (CRF)	Classif	Discrim	Param
<i>K</i> nearest neighbor classifier	Classif	Gen	Non
(Infinite) Mixture Discriminant analysis	Classif	Gen	Non
Classification and regression trees (CART)	Both	Discrim	Non
Boosted model	Both	Discrim	Non
Sparse kernelized lin/logreg (SKLR)	Both	Discrim	Non
Relevance vector machine (RVM)	Both	Discrim	Non
Support vector machine (SVM)	Both	Discrim	Non
Gaussian processes (GP)	Both	Discrim	Non
Smoothing splines	Regr	Discrim	Non

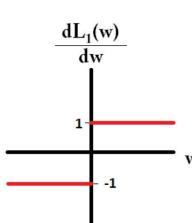
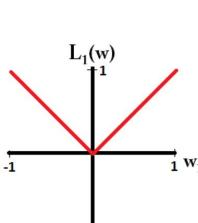


Compare and contrast L1/L2 regularization. What other names are they known by?

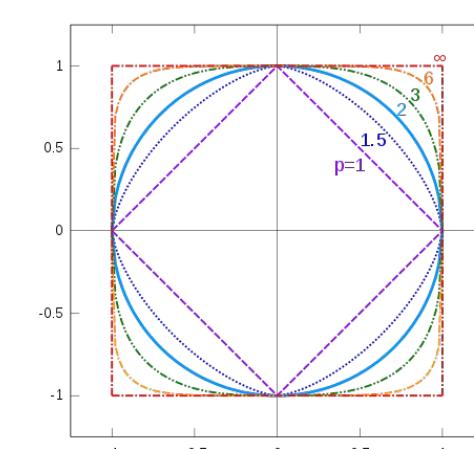
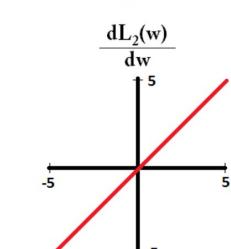
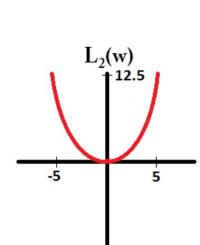
Suppose we have weights w_1, w_2, \dots, w_n . Then:

- L1 Regularization (i.e. Lasso regression):
 - Derivative/gradient is either 1, -1, or 0
 - Leads to sparsity (more weights being 0)
- L2 Regularization (i.e. Ridge regression):
 - Derivative/gradient is reduced linearly as the weight goes towards 0, so smaller gradient steps are taken (the motivation is diminished as you get closer)
 - Pushes weights closer to zero, but not sparsity

$$L_1(w) = \sum_i |w_i|$$



$$L_2(w) = \frac{1}{2} \sum_i w_i^2$$



Unit circles for different L_p norms.

What is data leakage and feature leakage?

- Leakage is the **use of information in the model training process which would not be expected to be available at deployment time**, causing the predictive scores (metrics) to overestimate the model's utility when run in a production environment
- **Data leakage** can be caused by improper sharing between the training and testing dataset. This can be caused by:
 - **Normalization bugs** (normalizing over entire dataset, not just training)
 - **Duplicated data** (e.g. if you oversample a minority class before partitioning)
 - **Non i.i.d. Data** (e.g. if you split grouped or time series data randomly)
- **Feature leakage** can be caused by redundant columns correlated with labels which by proxy, provide unfair signal to the actual label, and are not seen during deployment
 - Can also be caused by systemic training set bias, e.g. all picture of dogs are outdoors during daytime, but this is not the case in deployment
- Leakage can be detected if the model performs worse in production than in development. It should ideally be detected in shadow mode or canarying.

Statistical Data Processing

Define the following types of data processing: centering, normalization/feature scaling, decorrelation, standardization, whitening.

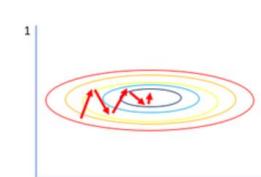
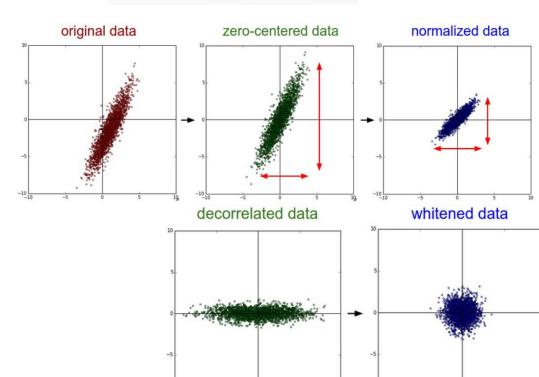
When is it necessary to standardize/normalize, and when is it not necessary?

- **Centering:** Subtract mean to each feature
- **Normalization/Feature Scaling:** Restrict all values in a dimension to be within a range [a,b].
- **Standardization:** Centers and sets the variances in the data to be 1 (diagonal of cov. Matrix becomes 1)
 - Can be done by dividing each dimension by its standard deviation
- **Decorrelation:** Removes the linear correlations in the data (cov. Matrix becomes diagonal)
 - Visually, rotates the data so the directions of most variation are aligned with the axes
 - PCA transformation does this
 - For regression, done only on the independent variables, not the dependent variables (otherwise, there would be nothing to regress)
- **Whitening:** Decorrelation + Standardization, data becomes uncorrelated and variances are 1 (identity cov. Matrix)
 - PCA, then dividing by the square root of the (now diagonal) covariance matrix does this
- Cases when scaling is **necessary**:
 - Helps gradient descent optimization **converge faster**
 - Prevents cost function contours from being too elliptical and more spherical
 - Needed since gradient descent is only 1st order, and does not use 2nd derivative based curvature information; descent directions are perpendicular to cost contour
 - PCA measures variance, so it's sensitive to the scale, if in different units
- Cases when scaling is **unnecessary**
 - Decision trees
 - When the scale conveys something important relative to other features

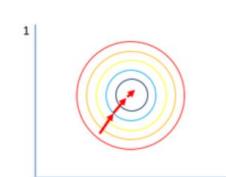
$$X' = a + \frac{(X - X_{\min})(b - a)}{X_{\max} - X_{\min}}$$

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

Feature Scaling for
[a,b] and [0,1]



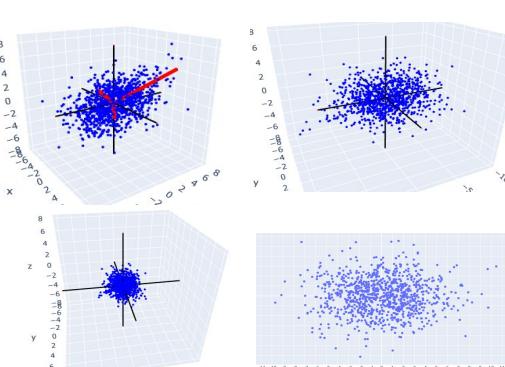
Gradient of larger parameter
dominates the update



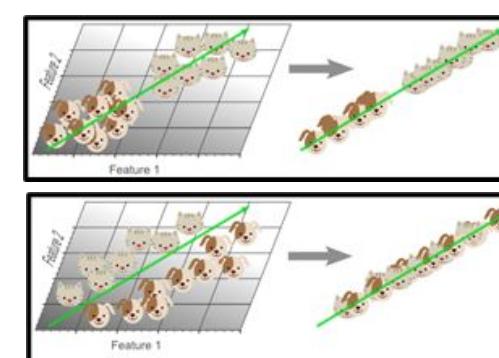
Both parameters can be
updated in equal proportions

Explain what PCA does, its steps and their first principles, what it's used for, and some pitfalls/relevant issues.

- PCA on a dataset gives an set of **orthogonal vectors** defining an orthogonal transformation (ie, “**rotation**”) that **transforms the data** such that scalar projecting the data onto the first axis of the transformation yields the greatest variance of any direction, projection onto the second axis yields the second greatest variance, etc.
- Using the set of axes (ie, “principal components”) that PCA provides:
 - The transformed dataset undergoes a change of basis, which renders the points **uncorrelated**
 - It can further be **whitened** by dividing each data point’s dimension by that dimension’s standard deviation. Then, the result is not only uncorrelated, but has **identity covariance**.
 - You can choose to only project the data to a few of the principal components, which allows for dimensionality reduction
- At a high level, the steps are:
 - Center the dataset
 - If features are of different units, standardize the data since PCA is sensitive to different ranges and variance
 - Gets insights of non-axis aligned covariance, beyond an ordering of the variables by their dimension-based variance
 - Intuitively, results should be same if you change from inches to cm and pounds to kg, etc
 - Compute the covariance matrix, and eigendecompose it so we have its eigenvectors (directions of greatest variance) and eigenvalues (magnitude of the variance)
 - Perform a change of basis using the top k eigenvectors (resulting in a $n \times k$ eigenvector matrix), to “rotate” the data points to be axes-aligned and uncorrelated
 - The dimensions not in the top k are “dropped”/ignored
 - You can also whiten it by (ie, multiplying with the covariance matrix to the $-\frac{1}{2}$ power)
- It’s nontrivial to decide how many principal components to keep; a heuristic is to plot r_k and pick a k such that sufficient variability is explained
- One assumption is that the most discriminative information is captured by the largest variance in the feature space. However, there are cases where the discriminative information actually resides in the directions of the smallest variance, such that PCA could greatly hurt classification performance. LDA avoids this.
- Theoretically, PCA is derived from trying to find directions which maximize the variance after projection. It turns out that eigenvectors solve this optimization problem.



Clockwise from upper left: points with eigenvectors plotted, decorrelated points, whitened points, projects projected to 2D



Top: Direction with most variance is discriminative

Bottom: Dimension with least variance is discriminative

$$r_k = \frac{\sum_{i=1}^k \sigma_i^2}{\sum_{i=1}^n \sigma_i^2}$$

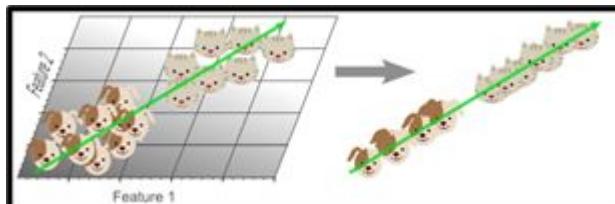
Explain what Linear Discriminant Analysis (LDA) is and what it's used for.

- Standard LDA finds the vector direction that maximizes the separability among 2 labeled categories
 - Contrast to PCA, which finds the directions which maximizes variance
 - LDA can be considered a “supervised version” of PCA
- Mathematically, we want to find a vector in the feature space that when the points are projected to it, the means of the two classes should be as **far apart as possible** and the **within-class variances small** as possible.
 - It can be shown that the solution to this optimization problem is w^*
- We can also generalize LDA to C many classes.
 - **Pairwise approach:** make $C(C,2)$ many classifiers/linear discriminants, and combine results by voting.
 - **One vs rest approach:** make C many “one vs rest” classifiers and combine results by voting.
- A generalized version of PCA for multiclass does provide a way to do dimensionality reduction for 1,..., C-1 dimensions
 - This is much less flexible than PCA which can vary the dimensionality from 1, ..., d

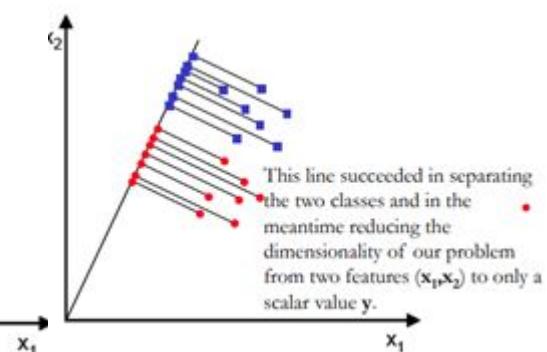
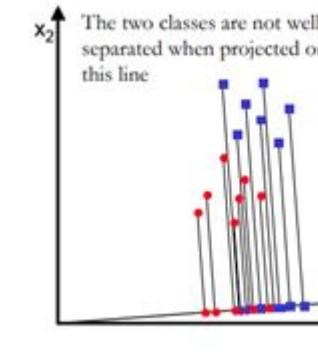
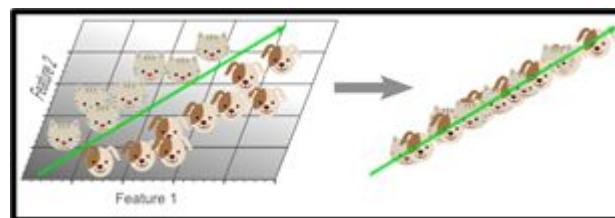
$$w^* = S_W^{-1}(\mu_1 - \mu_2)$$

$$S_W = S_1 + S_2, S_i = \sum_{j=1}^{N_i} (x_j - \mu_i)(x_j - \mu_i)^T$$

(S_i is the sample covariance without dividing by n-1)

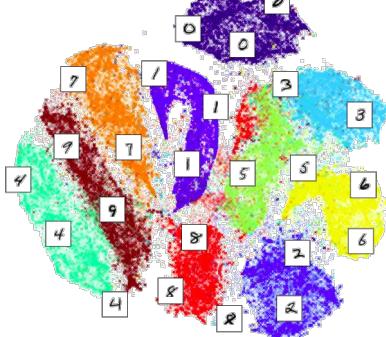


PCA would only work for the top, while LDA would work for both.



At a high level, how does t-SNE work?

- T-distributed stochastic neighbor embedding (t-SNE), proposed in 2008, is a variant of SNE (2002)
 - It allows for dimensionality reduction, with an emphasis on being able to visually preserve general notions of clusters, shape, and distance
 - Suppose you have points x_1, \dots, x_n . Then, the steps are:
 - Compute an $n \times n$ similarity matrix P . Each entry p_{ij} is a similarity score for x_i and x_j ; the sum of all p_{ij} is 1.
 - Similarity score is probability that x_i would pick x_j as its nearest neighbor, if picked in proportion to a gaussian centered at x_i
 - Initialize a point in the target dimension corresponding to each x_i . This can be random, or through PCA.
 - Until convergence:
 - Compute an $n \times n$ similarity matrix Q ; entry q_{ij} is a similarity score for the lower dimensional version of x_i and x_j , based on t-distribution distances. The sum of all q_{ij} is 1.
 - The fatter tails of the t-distribution help penalize distant pairs less heavily and prevent crowding in lower dimensions
 - Compute the KL divergence $KL(P||Q)$
 - Update the low dimensional points to minimize $KL(P||Q)$, using gradient descent
- Some rules of thumb:
 - Cluster sizes and distances between clusters after t-SNE should be taken with a grain of salt, and might not mean anything. The main focus is if there are clusters at all.



**t-SNE on MNIST, reducing
28*28=784 dimensional vectors
into 2 dimensions**

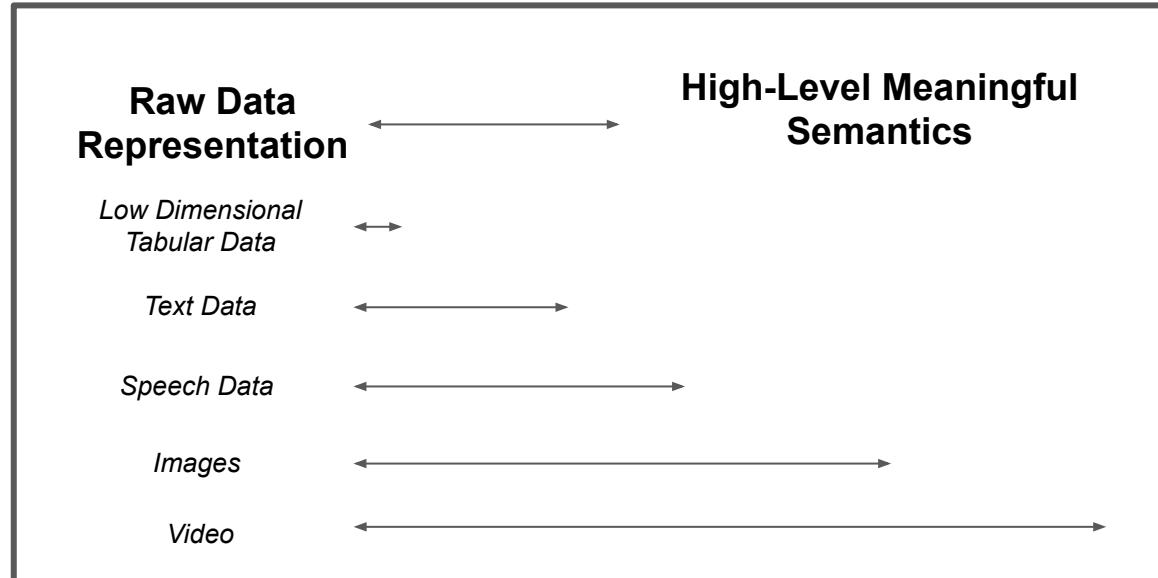
$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2 / 2\sigma^2)}$$

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

What does “representation” gap mean, at a high level?

- For different raw data modalities, there may be different sizes of gaps between it and a high-level meaning semantic needed for specific downstream tasks
 - This can be thought of as the semantic signal to noise ratio



Misc. Classical ML Models

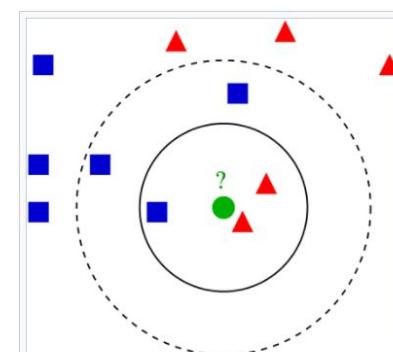
Give a list of classification algorithms.

- K-Nearest Neighbors (possibly after k=means clustering)
- Bayes classifier (e.g. naive Bayes)
- Neural Networks
- Decision trees (e.g. gradient boosted or random forest)
- SVM

How can K nearest neighbors be used for classification and regression? What are some other applications?

First, find the k nearest neighbors of your input in your training set. Then:

- For classification, output the class by taking a vote of the classes of those neighbors
 - A validation set can be evaluated by a confusion matrix and associated metrics like precision, recall, etc.
- For regression, average the values of those neighbors
 - A validation set can be evaluated by MSE error
- Applications include anomaly detection and retrieval (e.g. recommender systems)



Example of k -NN classification. The test sample (green dot) should be classified either to blue squares or to red triangles. If $k = 3$ (solid line circle) it is assigned to the red triangles because there are 2 triangles and only 1 square inside the inner circle. If $k = 5$ (dashed line circle) it is assigned to the blue squares (3 squares vs. 2 triangles inside the outer circle).

Explain the K-Means Clustering problem, an algorithm for it, and its applications.

Given a set of observations (x_1, x_2, \dots, x_n), where each observation is a d -dimensional real vector, k-means clustering aims to partition the n observations into k ($\leq n$) sets $S = \{S_1, S_2, \dots, S_k\}$ such that the within-cluster variances are minimized

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = \arg \min_{\mathbf{S}} \sum_{i=1}^k |S_i| \operatorname{Var} S_i$$

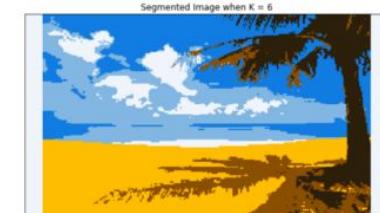
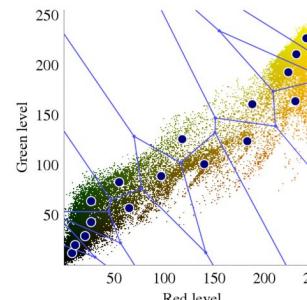
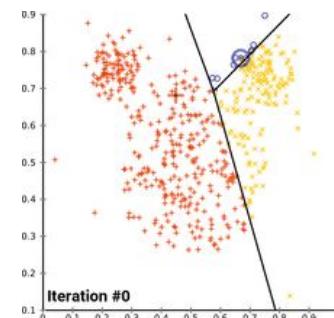
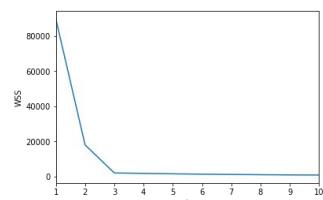
In practice, this is an NP-Hard problem and **Lloyd's algorithm** (which is not guaranteed to find the optimum) is usually used. Given an initial set of k means/centroids (eg, by randomly choosing from the observations), we repeat the following:

- **Assignment Step:** Assign each observation to the cluster with the nearest centroid (based on euclidean distance)
- **Update Step:** Recalculate the centroid by taking the mean of the points in each cluster

The value of k can be chosen based on “**the elbow method**”, which checks for when the **within-cluster-variance** suddenly drops for values of k .

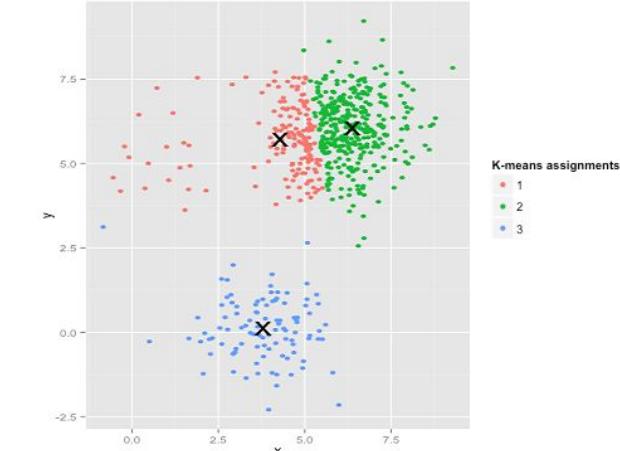
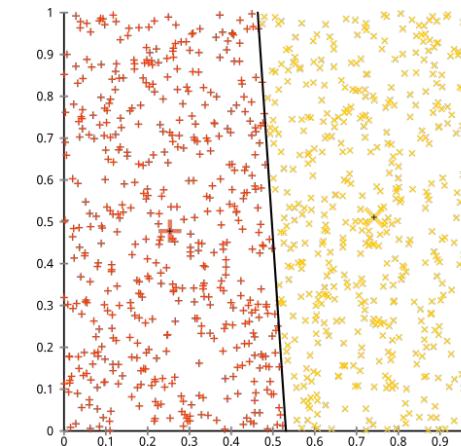
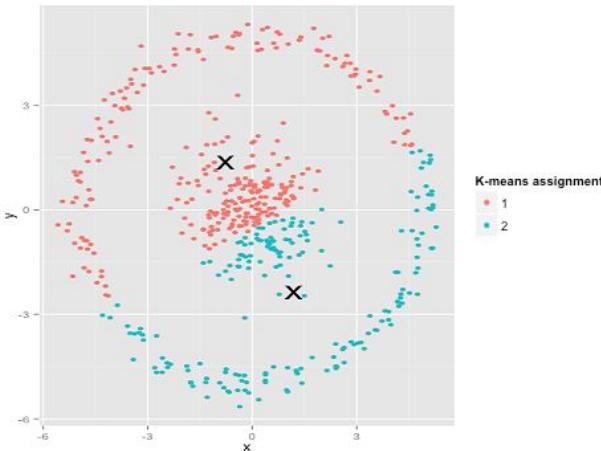
Applications to k-means:

- The resulting clusters can be used to discover/visualize groups within your dataset
 - Eg, discover types of customers, genres of books, fraud detection groups, groups of colors (for image segmentation)
- Voronoi diagram can be used to quantize/discretize the dataset's vectors into k bins
- Centroids can be used for k-NN classifier

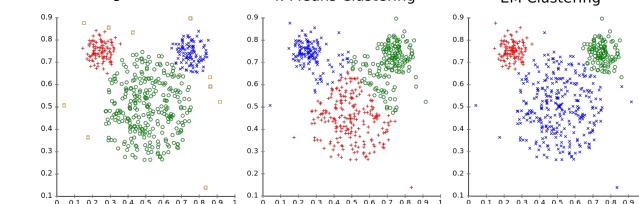


What are some failure cases of k-means?

- Bad initialization/optimization
- Spherical data (can be mitigated by embedding in a high dim space, or using polar coords)
- Data where there are no clusters
- Clusters of different variances and sparsity
 - More weight assigned to larger clusters
 - Instead use GMM, which generalizes k-means for non-homoscedastic and ellipsoid cases



Different cluster analysis results on "mouse" data set:
Original Data K-Means Clustering EM Clustering



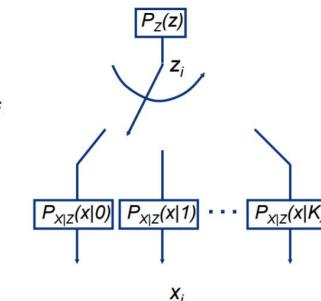
How are mixture models formulated, optimized, and used?

There are two types of random variables we are trying to model: the observed variable X and the hidden state variable Z we're "unaware of".

Observations are realized with a two step procedure:

1. z is sampled from Z , with pdf $P(Z=z)$
2. $x|z$ is sampled from $X|Z$, with pdf $P(X=x|Z=z)$.

$$P_X(x) = \sum_{c=1}^C P_x(x, Z=c) = \sum_{c=1}^C P_{X|Z}(x|c)P_Z(c) = \sum_{c=1}^C P_{X|Z}(x|c)\pi_c$$



The components of a mixture model are as follows:

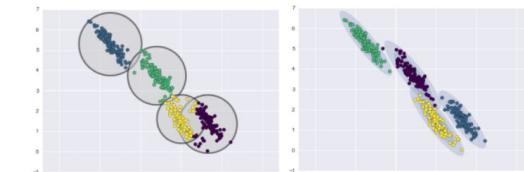
K	= number of mixture components
N	= number of observations
$\theta_{i=1\dots K}$	= parameter of distribution of observation associated with component i
$\phi_{i=1\dots K}$	= mixture weight, i.e., prior probability of a particular component i
ϕ	= K -dimensional vector composed of all the individual $\phi_{1\dots K}$; must sum to 1
$z_{i=1\dots N}$	= component of observation i
$x_{i=1\dots N}$	= observation i
$F(x \theta)$	= probability distribution of an observation, parametrized on θ
$z_{i=1\dots N}$	$\sim \text{Categorical}(\phi)$

Then, for example for a multivariate gaussian mixture model, the probability of an observation would be calculated as: $p(\theta) = \sum_{i=1}^K \phi_i \mathcal{N}(\mu_i, \Sigma_i)$

The EM algorithm is an iterative method to find the MLE or MAP, when models depends on unobserved latent variables. So, it's a natural fit for mixture models.

Applications of mixture models:

- Model a distribution for a dataset, and be able to understand the approximate probability of future values.
- Clustering the data into K groups (ie, more sophisticated way to do k-means with soft assignment; argmax for hard)
 - Given a data point, can compute the probability from each component distribution for each z by $P(X,Z)=P(X|Z)P(Z)$
 - It can be shown that k-means can fit spheres while GMMs can fit ellipsoids with non-identity covariances



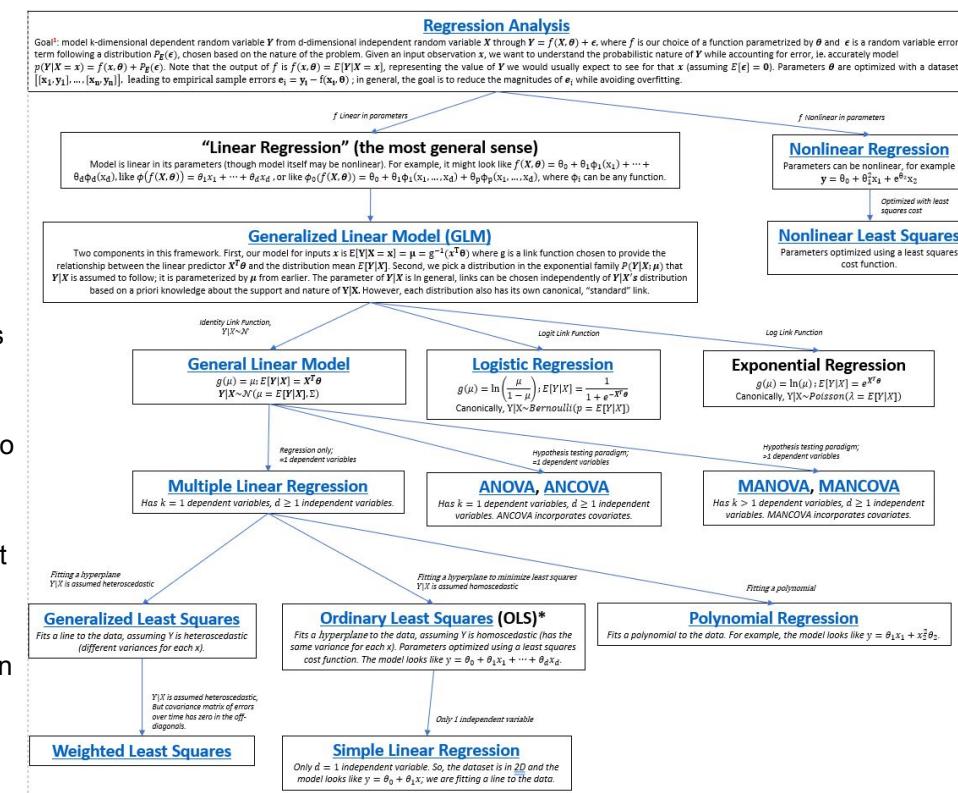
At a high level, describe the goal and hierarchy of regression analysis.

The goal of regression is to learn a mapping from independent variables X (i.e. features) to dependent variables Y (i.e. labels). The mapping has parameters, which are optimized by a training dataset of features their corresponding labels; we want this model to be generalizable, and avoid overfitting.

- **Nonlinear Regression:** model not linear in parameters, eg $y=x^*\theta^{\wedge}2$
- **“Linear Regression”, in the most general sense:** Model is linear in parameters (though, model itself can be nonlinear).

- **Generalized linear model**

- **Logistic regression:** For when the dependent variables are binary (follows a bernoulli distribution)
- **Multiple linear regression:**
 - **Ordinary Least Squares:** Fitting a hyperplane to data, assuming Y is homoscedastic. Optimized with least squares cost function.
 - **Generalized Least Squares:** Similar to OLS but allows Y to be heteroscedastic
 - **Polynomial Regression:** Fits a polynomial function to the data. (This is still linear regression since the estimator doesn't "know" x_i^2 is the square of x_i , it just thinks it's another variable. It's ultimately the same variable, we just transformed it to capture a curvilinear relationship. In general, we can transform individual independent variables arbitrarily, and even make multiple copies.)



Define, and compare/contrast logistic regression and softmax regression

Logistic Regression (Binary)

- Given an input vector X and vector of parameters β , we try to optimize the following using MLE:

$$P(Y = 1|X; \beta) = \frac{1}{1 + e^{-\beta^T X}}$$

$$P(Y = 0|X; \beta) = 1 - P(Y = 1|X; \beta)$$

- This is based off the sigmoid function $\frac{1}{1+e^{-x'}}$, which looks like this:

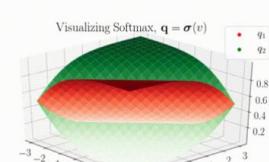


Multinomial Logistic Regression (i.e. "Softmax Regression")

- This generalizes binary logistic regression, and for K classes, each with a vector of parameters β_k , we optimize:

$$P(Y = c|X; \beta_1, \dots, \beta_K) = \frac{e^{\beta_c^T X}}{\sum_{k=1}^K e^{\beta_k^T X}}$$

- This is based off the softmax function, $\text{softmax}(k, x_1, \dots, x_n) = \frac{e^{x_k}}{\sum_{l=1}^n e^{x_l}}$, which looks like this:



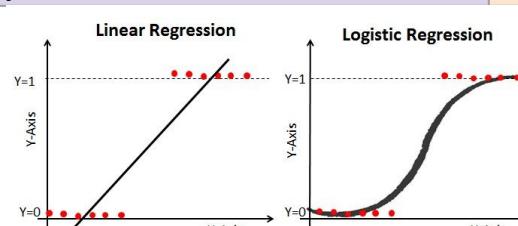
- Note that this is equivalent to binary logistic regression in the two-class case, with a change of parameters:

$$\frac{e^{\beta_1^T X}}{e^{\beta_1^T X} + e^{\beta_2^T X}} = \frac{1}{1 + e^{\beta_2^T X - \beta_1^T X}} = \frac{1}{1 + e^{(\beta_2 - \beta_1)^T X}}$$

- For multiclass, one can also do K binary one-vs-all logistic regressions
 - However, you would not be able to compare probabilities of two classes
 - There are some theoretical advantages when the model is fit simultaneously, such as different goodness-of-fit tests and informative residuals. However, there are also some advantages when they are fit separately, and provide information at the individual-case level.

Compare/contrast ordinary least squares linear regression and logistic regression.

	Ordinary Least Squares Linear Regression	Logistic Regression
High Level Purpose	Fits a hyperplane to the data, assuming the dependent variable is homoscedastic	Fits a sigmoid to the data where the dependent variable is binary and follows a bernoulli distribution..
Model (theta is a vector of parameters)	$y = X^T \theta = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$	$h_{\theta}(X) = \frac{1}{1 + e^{-\theta^T X}} = \Pr(Y = 1 X; \theta)$ $\Pr(Y = 0 X; \theta) = 1 - h_{\theta}(X)$
Loss function (recall, MLE has many good theoretical properties to ensure asymptotic optimality, eg unbiasedness)	MSE (ie “least squares”), which can be shown to be equivalent to the MLE. $S(b) = \sum_{i=1}^n (y_i - x_i^T b)^2$	MLE, which is the binary cross entropy: $-\frac{1}{N} \sum_{i=1}^N \log \Pr(y_i x_i; \theta)$ $\frac{1}{N} \sum_i [-(y_i) \log(\Pr(y_i x_i; \theta)) - (1 - y_i) \log(1 - \Pr(y_i x_i; \theta))]$ In principle one can use MSE, but that loses some of the properties of MLE.
Optimization Procedure	Convex, so can be optimized in closed form, but sometimes a gradient descent approximation is used for computational efficiency	No closed form; need gradient descent. So, scaling/normalization often helps.



Explain what Support Vector Machines (SVMs) are, how they work, and what their pros/cons are.

SVMs are **non-probabilistic binary linear classifiers**, which solves an optimization problem to find a hyperplane to linearly separate labeled training data with maximum margin.

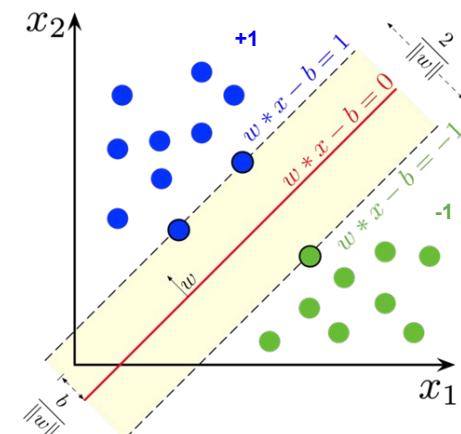
To learn an SVM classifier:

- Standardize the data
- There are 3 hyperplanes of form $wx - b$, as shown in the figure; their distance is $2/\|w\|$. We want to maximize that distance while maintaining that data points lie on their respective side of the margin as much as possible. To do this, the following is minimized:

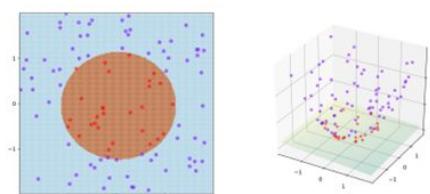
$$\lambda\|w\|^2 + \left[\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i - b)) \right]$$

This follows since $1 - y_i(w^T x_i - b) \leq 0$ should hold if a point is on the correct side.

- Lambda determines the tradeoff between ensuring the data points lie on their correct side (may lead to overfitting) and large margin (may lead to underfitting/ignoring features). Note that if all points are perfectly linearly separable, then the second term should be zero.
- The resulting w^* , b^* (e.g. from SGD or quadratic programming) can be used for classification: $x \mapsto \text{sign}(w^{*T} x - b^*)$



Additionally, the kernel method can be used by replacing $w^T x_i = \langle w, x_i \rangle$ with $K(w, x_i) = \langle \Phi(w), \Phi(x_i) \rangle$



Some strengths of SVMs:

- Makes little to no assumptions about data distribution
- Training is easy; no local optima since the function is **convex**
- Effectiveness scales relatively well to high dimensions
- Relatively **robust to outliers**, due to maximal margins & because non-support data points are effectively “ignored”
- **Fast** at inference

Some weaknesses:

- **Requires a good kernel function**
- **No probability output for estimate**

Compare/Contrast SVMs and LDA

SVM:

- Makes **no assumptions** about the data, making it flexible
- Optimized over a subset of the data: only data points that lie on the separating margin (“**support vectors**”)
- Able to use **kernels** to transform into nonlinear classifier

LDA:

- Optimality requires that data is **normally distributed**
- Makes **use of the entire data set**
- Can be used for **feature selection/dimensionality reduction**

What is Bayes' classifier? What is the Naive Bayes classifier? What are their pros and cons?

Bayes' classifier takes a probabilistic approach, centered around Bayes' rule:

$$p(Y|X) = \frac{p(X|Y) p(Y)}{p(X)} \Leftrightarrow \text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

Ultimately, given some x we want to choose the most likely y -label. The decision rule is:

$$\hat{y} = \operatorname{argmax}_{i \in \{1, \dots, k\}} p(x|Y=i) * p(Y=i)$$

Note that the denominator is ignored for classification purposes, since it's just a constant. To solve this, one needs to estimate the distribution for $X|Y$ and Y .

- Y follows a multinoulli distribution and can be estimated with MLE, or assumed to be uniform across all classes.
- $X|Y$ is a n-dimensional distribution, and needs to be fit. For example, one option is to fit k -many (one for each class conditional) n-dimensional gaussians, and estimate the parameters with MLE.

Naïve Bayes

If we assume that the X_1, \dots, X_n in X are independent from one another, then fitting $X|Y$ becomes easier. In this case,

$$\begin{aligned} p(X|Y)p(Y) &= P(X, Y) = p(X_1, \dots, X_n, Y) = p(X_1|X_2, \dots, X_n, Y)p(X_2, \dots, X_n, Y) \\ &= p(X_1|X_2, \dots, X_n, Y) \cdots p(X_{n-1}|X_n, Y)p(X_n|Y)p(Y) \\ &= p(Y) \prod_i^n p(X_i|Y), \text{ due to independence} \end{aligned}$$

This makes estimating the distribution of $X|Y$ easier and less data-dependent.

For example, we can fit $k * n$ many 1-dimensional gaussians and just multiply the probabilities, which avoids the curse of dimensionality.

Pros of Bayes' Classifier:

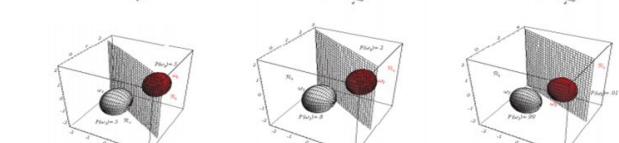
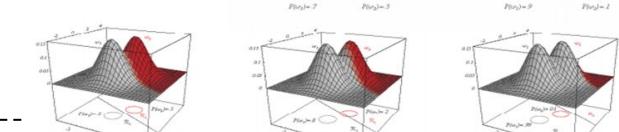
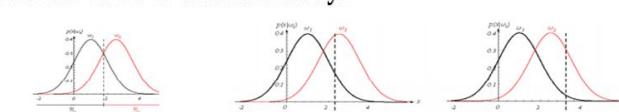
- Technically optimal, if distribution estimates are accurate
- Generative and gives "exact" per-class probabilities (e.g. unlike SVMs)

Cons:

- Hard to estimate probability distribution for $X|Y$ (need a lot of data)
- Naive bayes assumption requires independence, which probably won't hold in most cases

Suppose you use gaussians to model $X|Y$. The figure shows the binary classification case for $x \in R^1, R^2, R^3$ in the first, second, and third rows respectively.

The decision boundary given by the BDR is a hyperplane, and is weighed by the distribution of Y . As the black class becomes more likely (left to right), it occupies more space and pushes the decision boundary further away from it.



Give an example of Bayes classifier for sentiment analysis on the following dataset, using laplace smoothing, and a bag-of-words model.
(unique =20)

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no fun

$$\hat{\theta}_i = \frac{x_i + \alpha}{N + \alpha d} \quad (i = 1, \dots, d)$$

* note: denominator is necessary since you add 1 to all d "buckets" or categories from the multinomial distribution

$$P(-) = \frac{3}{5} \quad P(+) = \frac{2}{5} \quad \begin{aligned} P(\text{"predictable"}|-) &= \frac{1+1}{14+20} & P(\text{"predictable"}|+) &= \frac{0+1}{9+20} \\ P(\text{"no"}|-) &= \frac{1+1}{14+20} & P(\text{"no"}|+) &= \frac{0+1}{9+20} \\ P(\text{"fun"}|-) &= \frac{0+1}{14+20} & P(\text{"fun"}|+) &= \frac{1+1}{9+20} \end{aligned}$$

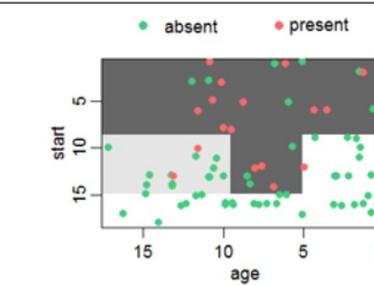
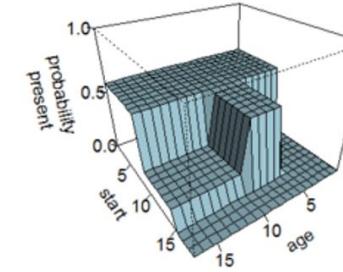
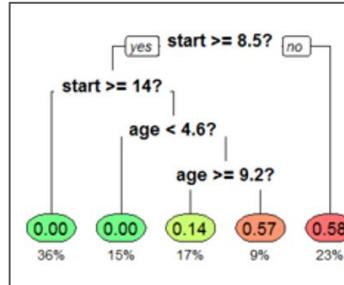
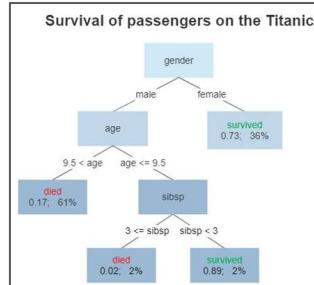
$$P(-)P(S|-) = \frac{3}{5} \times \frac{2 \times 2 \times 1}{34^3} = 6.1 \times 10^{-5}$$

$$P(+)P(S|+) = \frac{2}{5} \times \frac{1 \times 1 \times 2}{29^3} = 3.2 \times 10^{-5}$$

The model thus predicts the class *negative* for the test sentence.

What are decision trees? Explain its pros/cons.

Decision trees are binary trees where each non-leaf node directs the input vector to either its left or right node, based on two per-node parameters: the feature dimension, and threshold c . The leaf nodes are classification or regression prediction outputs.



Decision trees can be learned using a training dataset $\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^K$.

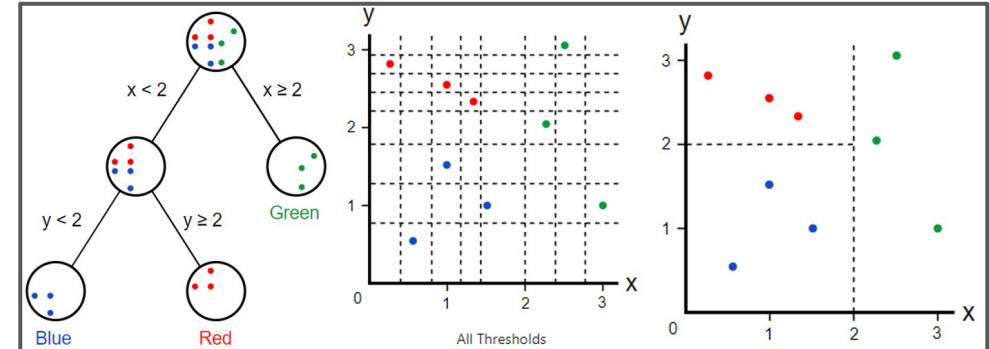
- Recursively, decision nodes are added to the tree, where each node is parametrized by the feature dimension in $\{1, \dots, k\}$, and a threshold c .
 - There are several ways to decide this, but a common approach is just brute force: try every possible threshold for every dimension, and pick the one with the highest **Gini Gain**. At a high level, this allows us to pick the “best” split.
- A node is no longer split if all points in it have the same class, or if all splits are equally good.

Pros:

- Decision Trees are easily **interpretable**, and mirrors human decision making
- **Built-in feature selection**: unused features can be discarded, and features on top are the most informative
- Handles both categorical and numerical data easily; robust to scale & heterogeneous features
- Capable of learning **nonlinear**, complicated decision boundaries
- **Simple** to use “out-of-the-box”

Cons:

- Can be very prone to **overfitting**
- Learning an optimal decision is known to be NP-complete, so in practice all algorithms will be **suboptimal** (e.g. using greedy methods)



What is a Random Forest?

Random Forests are essentially ensembles of T decision trees. This generally reduces variance/overfitting.

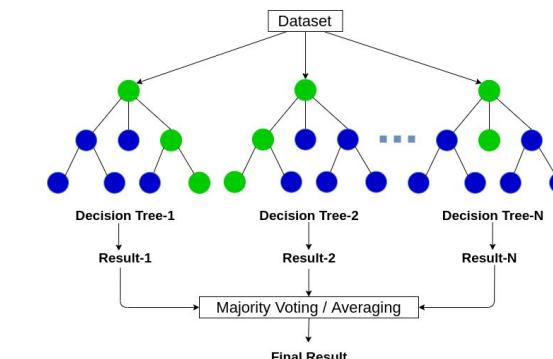
Given a training dataset of size N, for each decision tree:

- Sample, **with replacement**, N training examples
- Train a decision tree with those N samples
 - In the training, each time we make a new decision node, only try a subset of the features (e.g.), to make each tree more random and reduce correlation between trees.

Sampling with replacement and using only a subset of features helps inject more randomness. Then, the final result is either averaged (for regression) or voted upon (for classification).

Some terminology:

- Sampling with replacement is called **bootstrapping**
- The technique of creating an ensemble using bootstrapping is called **bootstrap aggregating (bagging)**
- The technique of only trying a subset of the features is called feature bagging.



Explain Adaboost, and at a high level, gradient boosting. What are its pros and cons?

- Boosting is a meta-algorithm to learn strong classifiers by combining an **ensemble** of simpler, “weak” ones.
 - The weak ones do not need to perform very well; only at least better than random chance. A popular choice are **decision stumps** (decision trees with 3 nodes).
 - The best weak learners are iteratively added with a weight into the overall classifier, and counts as a “vote” towards the overall classification. Another set of weights are also used for the training data points, to encourage weak learners to address misclassified points

Adaboost pseudocode:

Suppose you have training data $\{(x_i, y_i)\}_{i=1}^N$, $x_i \in \mathbb{R}^K$, $y_i \in \{0,1\}$. Additionally, we have a set of weak classifiers $F = \{f_m(x)\}$, $f_m(x) \in \{-1,1\}$.

Then, the Adaboost classifier is

$$g(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t f_t(x)\right).$$

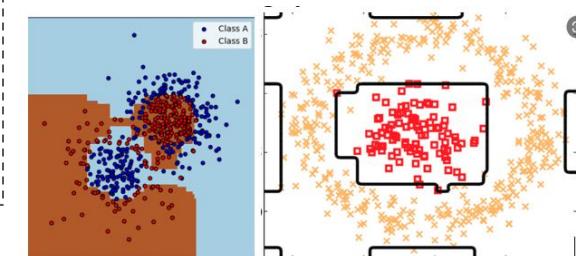
Here, the α_t are learned as follows, for T iterations:

- Assign an initial weight $w_i^1 = 1/N$ for each datapoint x_i
- For iteration t=1 to T:
 - Find a weak learner f_t such that the normalized, weighted sum error for misclassified points is minimized:

$$\epsilon_m = \frac{\sum_{i=1}^N w_i^{(m)} I(f_m(x_i) \neq y_i)}{\sum_i w_i^{(m)}}, \quad I(f_m(x_i) \neq y_i) = 1 \text{ if } f_m(x_i) \neq y_i \text{ and } 0 \text{ otherwise}$$

$$\alpha_t = \ln\left(\frac{1-\epsilon_m}{\epsilon_m}\right)$$

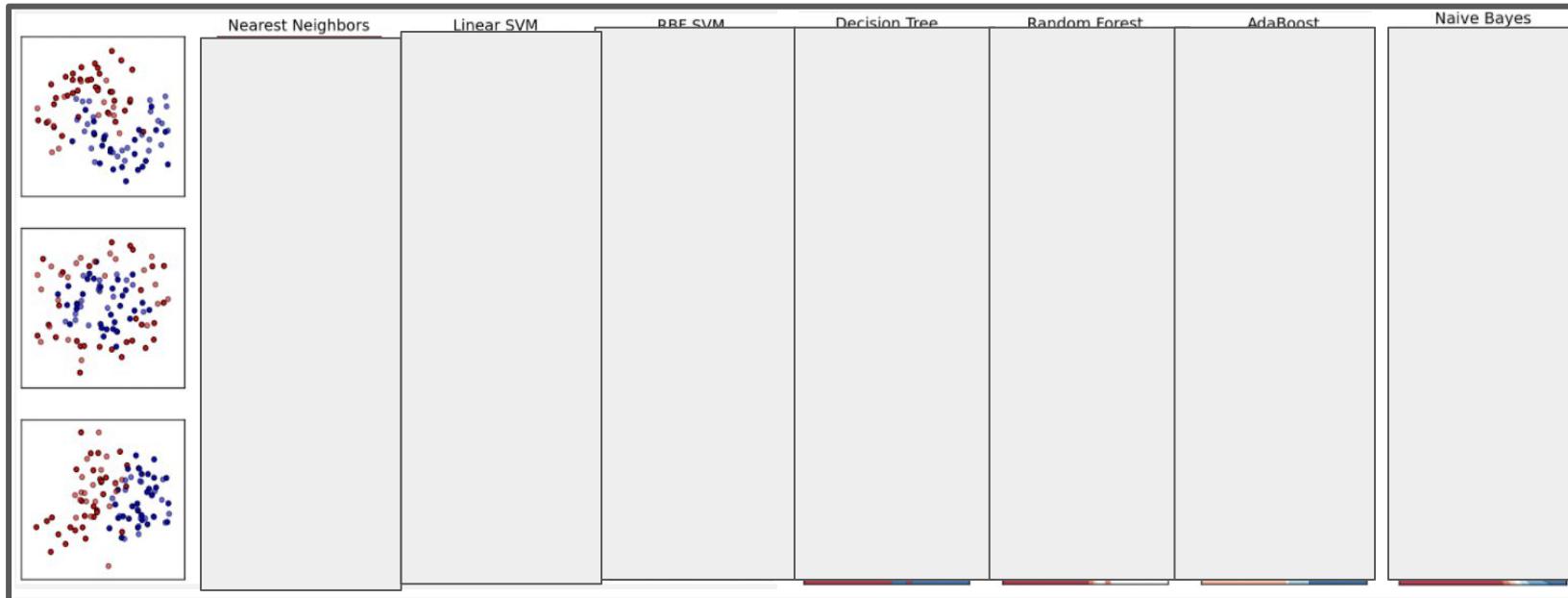
$$\text{Reassign all weights as } w_i^{t+1} = w_i^t e^{\alpha_m I(f_m(x_i) \neq y_i)}$$



- **Pros of Adaboost:**
 - It can work well on very small datasets, eg those under 100 samples.
 - Adaboost (with decision trees as the weak learners) is often referred to as the **best out-of-the-box classifier** (according to Wikipedia)
 - It can learn highly non-linear decision boundaries
- **Cons of Adaboost:**
 - Harder to tune than random forests
 - Easier to overfit than random forests
- Adaboost can be interpreted as gradient descent in function space, where at each descent iteration, the alphas are the step sizes and the chosen weak learners are the gradient.
- Adaboost is a special case of **Gradient Boosting**, where Adaboost uses a special exponential loss function. Gradient Boosting can handle a variety of loss functions. Note that Gradient Boosting was developed after Adaboost.

For the 3 given datasets, visually the decision boundaries for: NN, linear SVM, RBF SVM, Decision Trees, Random Forest, AdaBoost, and Naive Bayes.

In relative terms, which ones do you think would perform the best?



At a high level, what is a Bayes Filter? What are Kalman filters?

- **Bayes Filters** estimate an unknown probability density function (PDF) recursively over time, using incoming measurements and a mathematical process model.
 - For noisy measurements are observed over time, Bayes filters give estimates that tend to be more accurate than with a single measurement alone, using temporal smoothing
 - They allow one to integrate a mathematical motion model, prior knowledge, and noisy measurements into a predictive framework which accounts for noise
 - Bayes filters take the Markov assumption: the current robot/environment state only depends on the previous one, and not any more timesteps back.
- **Kalman filters** are an important type of Bayes Filter which takes several assumptions:
 - The prior pdf is Gaussian
 - The motion & observation models are linear in the state, and affected by Gaussian noise
 - The process noise and measurement noise are independent of each other, of the state, and across time
- Given the previous state, current control input, and current measurements, there are two steps which are repeated:
 - **Prediction Step:** Using the previous state, motion model, and noise model, the filter produces a prediction of what the current state is.
 - **Update Step:** Using the next (noisy) measurement, the current estimated state is refined/updated.

Predict:

$$\hat{x}_{t|t-1} = F_t \hat{x}_{t-1|t-1} + B_t u_t$$

$$P_{t|t-1} = F_t P_{t-1|t-1} F_t^T + Q_t$$

Update:

$$\hat{x}_{t|t} = \hat{x}_{t|t-1} + K_t (y_t - H_t \hat{x}_{t|t-1})$$

$$K_t = P_{t|t-1} H_t^T (H_t P_{t|t-1} H_t^T + R_t)^{-1}$$

$$P_{t|t} = (I - K_t H_t) P_{t|t-1}$$

Current state updated using previous state, motion models, and control input

Current state covariance increases, using estimation error and process error distributions

Refine current state using observations; this decreases covariance.

The Kalman gain intuitively describes how much we trust our observations.

It is multiplied with the "innovation term" $y - Hx$, which is the difference between the observation and prediction.

\hat{x} : Estimated state.

F : State transition matrix (i.e., transition between states).

u : Control variables.

B : Control matrix (i.e., mapping control to state variables).

P : State variance matrix (i.e., error of estimation).

Q : Process variance matrix (i.e., error due to process).

y : Measurement variables.

H : Measurement matrix (i.e., mapping measurements onto state).

K : Kalman gain.

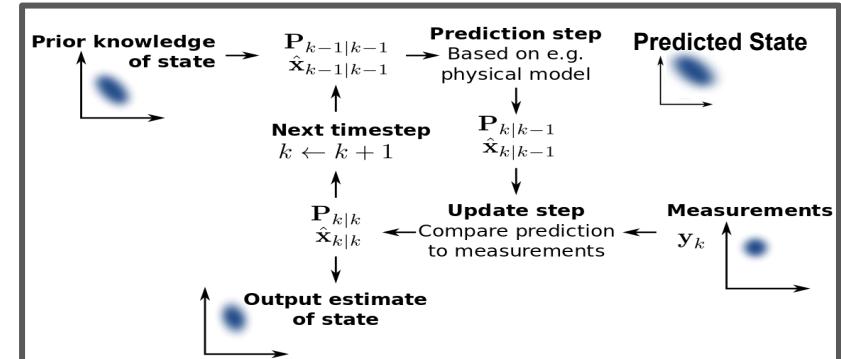
R : Measurement variance matrix (i.e., error from measurements).

- **Pros of Bayes filters:**

- Well understood, very general formulation with many applications for time-series analysis in engineering, robotics (motion planning/control), signal processing, finance, and economics
- In CV, can be used for temporal smoothing, eg depth estimation or tracking
- Can predict future values (using multiple predict steps)
- Provides a built-in way to integrate prior knowledge about state and motion model
- Under assumptions, can be proven that it's the best possible linear estimator in the minimum mean-square-error sense
- Can be run in real time

- **Cons of Bayes filters:**

- Requires motion model (e.g. based on physics), which may not be apparent. In these cases, an RNN/LSTM might work better. In fact, they can substitute for a Kalman filter in theory, and are more general/can be arbitrarily complex with many parameters.
 - An alternative is to keep the motion model as identity.
- Assumptions might not hold



Question goes here

Answer goes here