

SVMs, Kernelized SVMs

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Kai-Wei Chang, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Outline

- 1 Review of last lecture
- 2 SVM – Hinge loss
- 3 Kernelized SVMs
- 4 Evaluating ML algorithms

Support Vector Machine

- A linear classifier (hyperplane) that maximizes the margin.
- An alternative view of SVMs.

Outline

- 1 Review of last lecture
- 2 SVM – Hinge loss**
- 3 Kernelized SVMs
- 4 Evaluating ML algorithms

A general view of supervised learning

Definition Assume $y \in \{-1, 1\}$ and the decision rule is $h(\mathbf{x}) = \text{SIGN}(\underline{a(\mathbf{x})})$ with $\underline{a(\mathbf{x})} = \underline{\mathbf{w}^T \phi(\mathbf{x}) + b}$

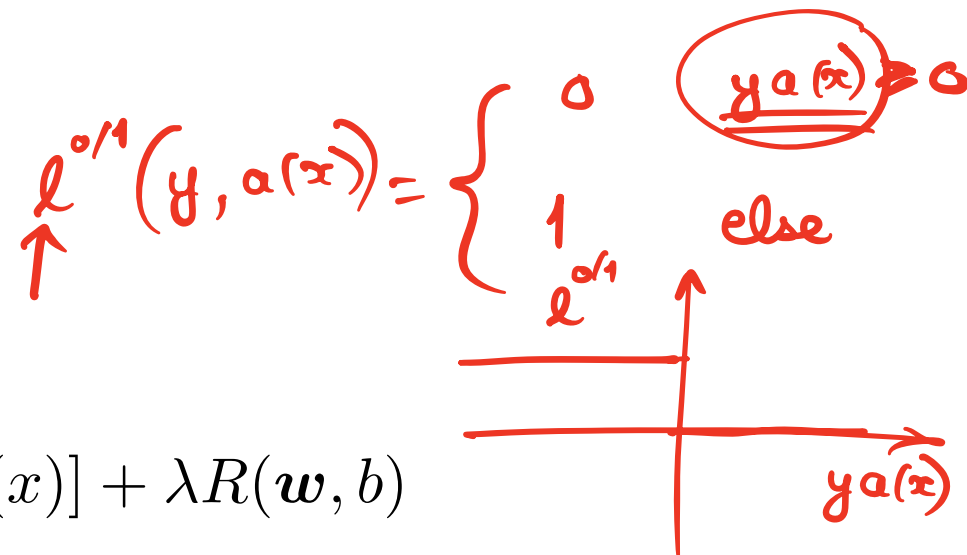
For classification: 0/1 loss

$$\underline{\ell^{0/1}(y, \underline{a(\mathbf{x})})} = \begin{cases} 0 & \text{if } \underline{ya(\mathbf{x})} \geq 0 \\ 1 & \text{otherwise} \end{cases}$$

Minimize weighted sum of empirical risk and regularizer

$$\begin{aligned} & \arg \min_{\mathbf{w}, b} \left[\underline{R^{\text{EMP}}[h_{\mathbf{w}, b}(x)]} + \underline{\lambda R(\mathbf{w}, b)} \right] \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \underline{\ell(y_n, a(\mathbf{x}_n))} + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

Minimize weighted sum of empirical risk and regularizer

$$\ell^{0/1}(y, a(x)) = \begin{cases} 0 & \text{if } y = a(x) \\ 1 & \text{else} \end{cases}$$


$$\begin{aligned} & \arg \min_{\mathbf{w}, b} R^{\text{EMP}}[h_{\mathbf{w}, b}(x)] + \lambda R(\mathbf{w}, b) \\ &= \arg \min_{\mathbf{w}, b} \frac{1}{N} \sum_n \ell(y_n, a(\mathbf{x}_n)) + \lambda R(\mathbf{w}, b) \end{aligned} \quad (1)$$

- Problem with minimizing the 0/1 loss ?

Hinge loss

Definition Assume $y \in \{-1, 1\}$ and the decision rule is $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$ with $a(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$,

$$\ell^{\text{HINGE}}(\underline{y}, \underline{a(\mathbf{x})}) = \begin{cases} 0 & \text{if } \underline{ya(\mathbf{x})} \geq \underline{1} \\ \underbrace{1 - \underline{ya(\mathbf{x})}} & \text{otherwise} \end{cases}$$

Intuition

$$\underline{ya(\mathbf{x})} < 0$$

$$ya(\mathbf{x}) \geq 0$$

$$ya(\mathbf{x}) < 1$$

Hinge loss

Definition Assume $y \in \{-1, 1\}$ and the decision rule is $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$ with $a(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$,

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \begin{cases} 0 & \text{if } ya(\mathbf{x}) \geq 1 \\ \underbrace{1 - ya(\mathbf{x})}_{ya(\mathbf{x}) < 0} & \text{otherwise} \end{cases}$$

Intuition

- No penalty if raw output, $a(\mathbf{x})$, has same sign and is far enough from decision boundary (i.e., if 'margin' is large enough)
- Otherwise pay a growing penalty, between 0 and 1 if signs match, and greater than one otherwise

Hinge loss

Definition Assume $y \in \{-1, 1\}$ and the decision rule is $h(\mathbf{x}) = \text{SIGN}(a(\mathbf{x}))$ with $a(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$,

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \begin{cases} 0 & \text{if } ya(\mathbf{x}) \geq 1 \\ 1 - ya(\mathbf{x}) & \text{otherwise} \end{cases}$$

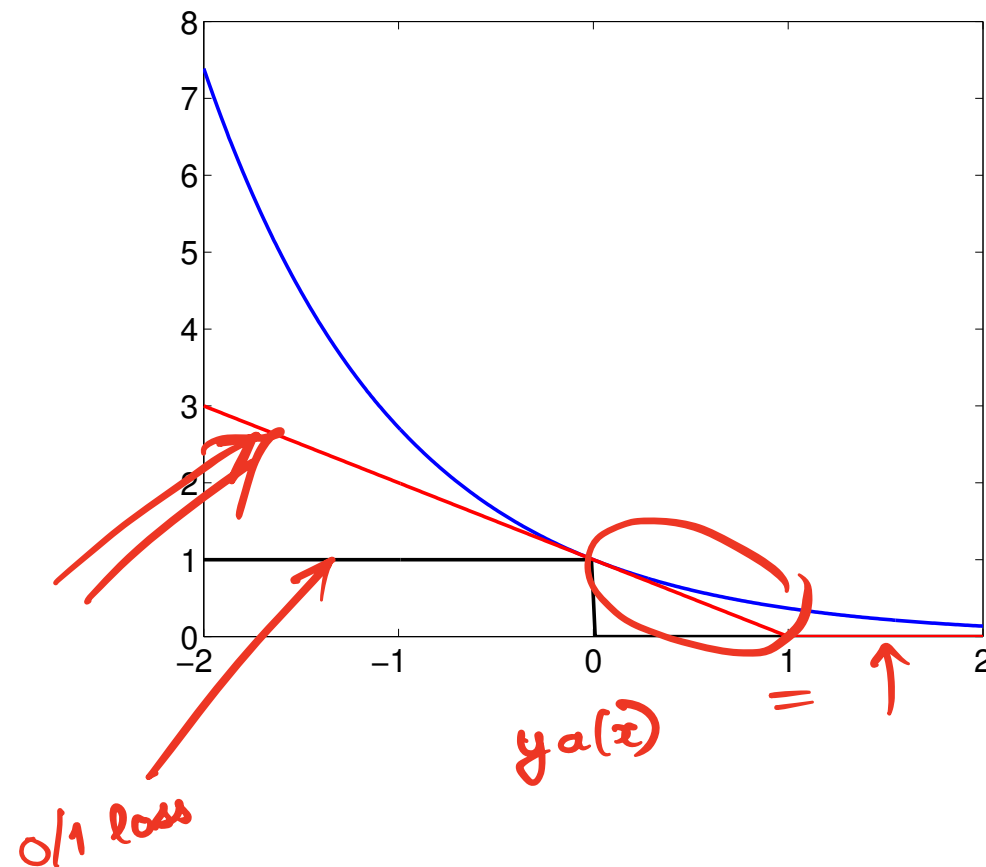
Intuition

- No penalty if raw output, $a(\mathbf{x})$, has same sign and is far enough from decision boundary (i.e., if ‘margin’ is large enough)
- Otherwise pay a growing penalty, between 0 and 1 if signs match, and greater than one otherwise

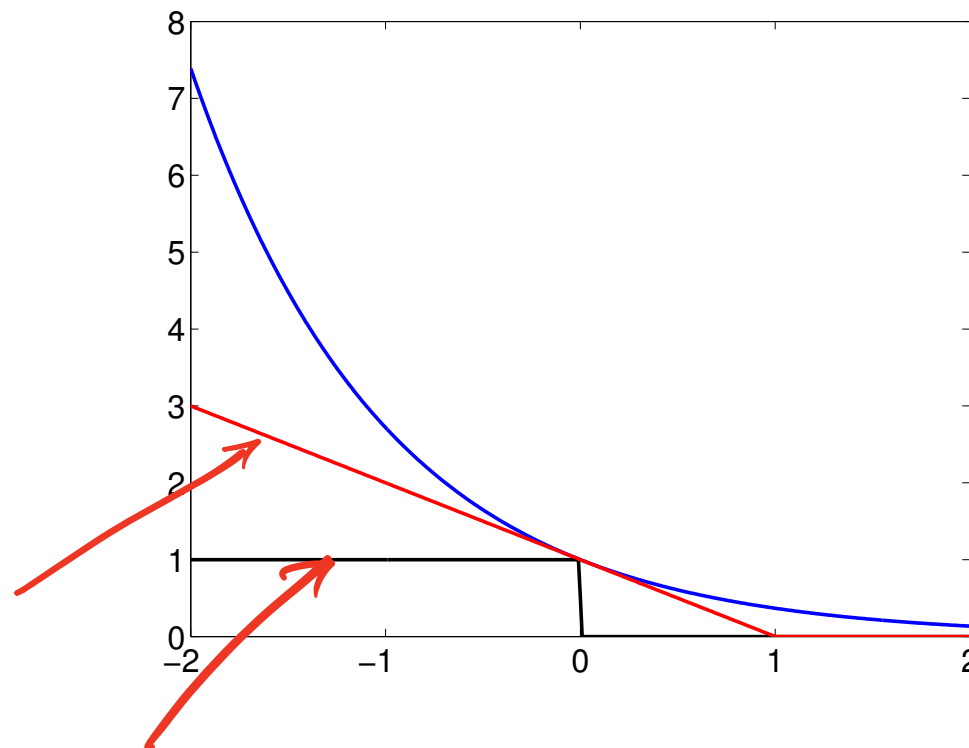
Convenient shorthand

$$\ell^{\text{HINGE}}(y, a(\mathbf{x})) = \max(0, 1 - ya(\mathbf{x})) = (1 - ya(\mathbf{x}))_+$$

Visualization and Properties

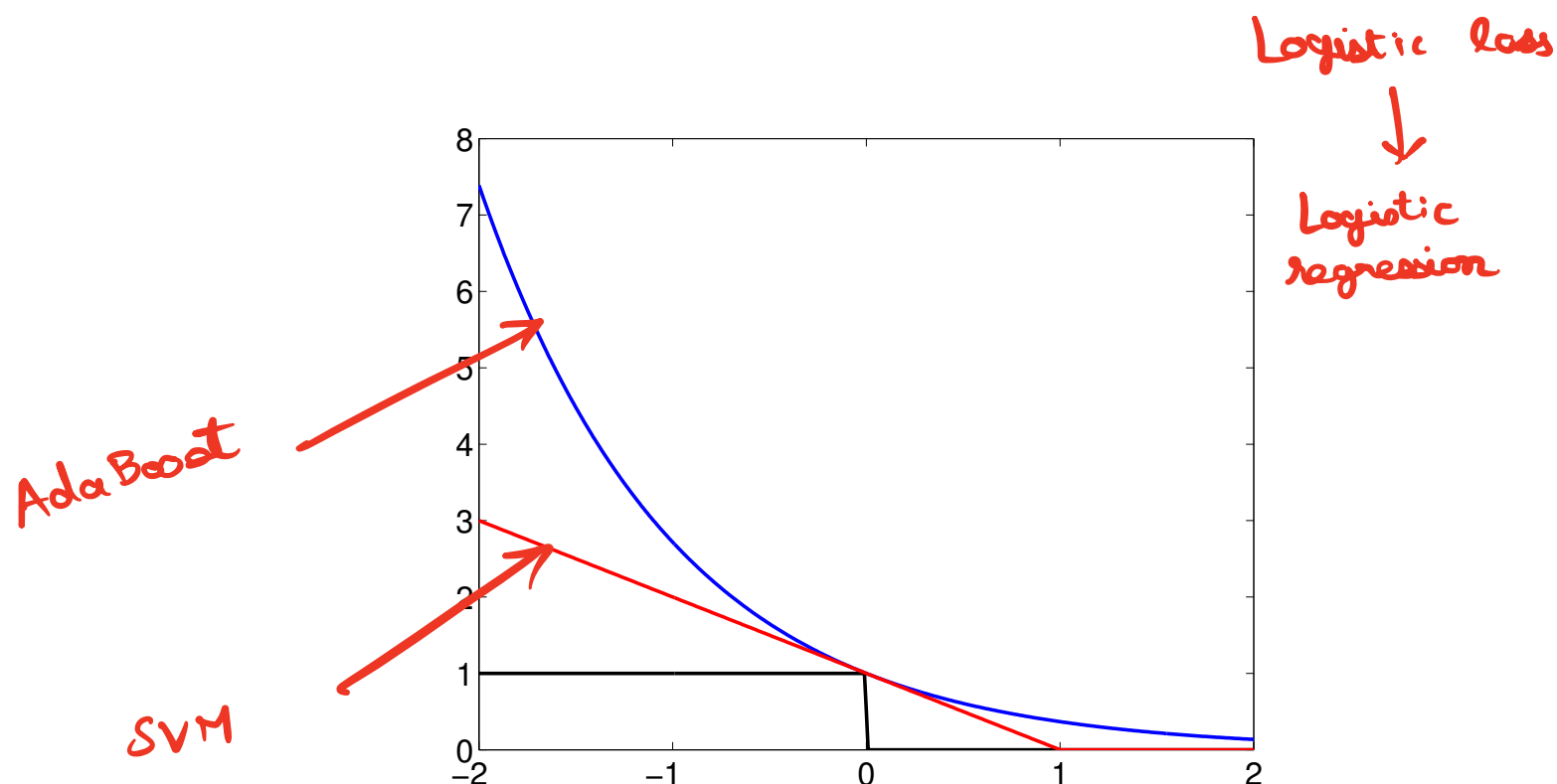


Visualization and Properties



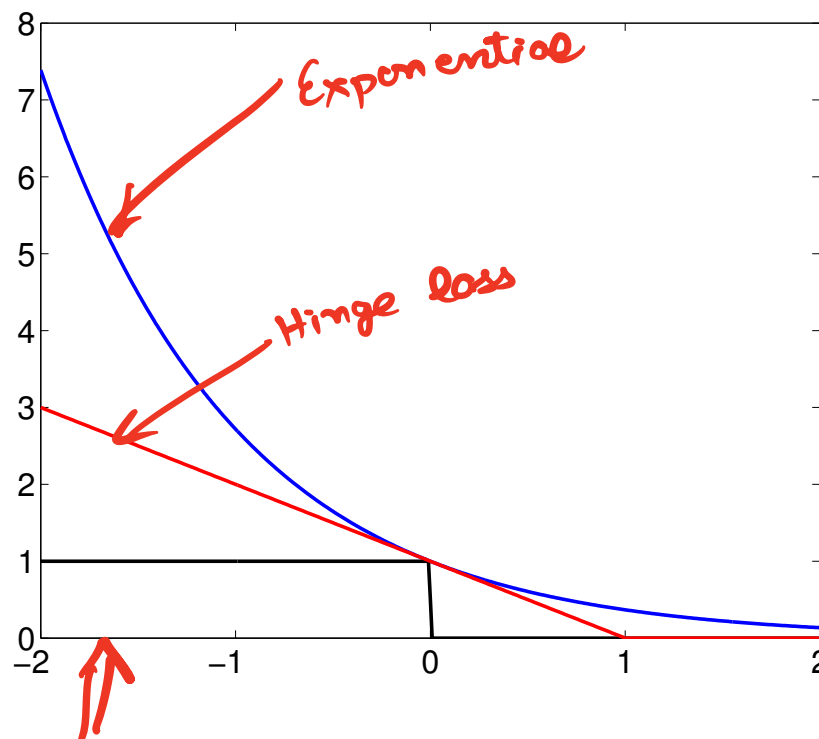
- Upper-bound for 0/1 loss function (black line)
- We use hinge loss as a *surrogate* to 0/1 loss – Why?

Visualization and Properties



- Upper-bound for 0/1 loss function (black line)
- We use hinge loss as a *surrogate* to 0/1 loss – Why?
- Hinge loss is convex, and thus easier to work with.

Visualization and Properties



- Other **surrogate losses** can be used, e.g., exponential loss for Adaboost (in blue), logistic loss (not shown) for logistic regression
- Hinge loss less sensitive to outliers than exponential (or logistic) loss

Primal formulation of support vector machines (SVM)

Minimizing the total hinge loss on all the training data with l_2 regularization

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Analogous to l_2 regularized least squares, as we balance between two terms (the loss and the regularizer).

Hinge loss

(\mathbf{x}_n, y_n)

$$\max(0, 1 - y_n a(\mathbf{x}))$$
$$a(\mathbf{x}) = [\mathbf{w}^T \phi(\mathbf{x}) + b]$$

Primal formulation of support vector machines (SVM)

Minimizing the total hinge loss on all the training data with l_2 regularization

$$\Rightarrow \min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Analogous to l_2 regularized least squares, as we balance between two terms (the loss and the regularizer).

Previously, we used geometric arguments to derive:

$$\begin{aligned} & \min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ & \text{s.t.} \quad y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \text{ and } \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

Do these yield the same solution?

Recovering our previous SVM formulation

Minimizing the total hinge loss on all the training data

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Define $C = 1/\lambda$:

$$\min_{\mathbf{w}, b} \underline{C} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

Recovering our previous SVM formulation

Minimizing the total hinge loss on all the training data

$$\min_{\mathbf{w}, b} \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Define $C = 1/\lambda$:

$$\min_{\mathbf{w}, b} C \sum_n \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

Define ξ_n = $\max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b])$

$$\min_{\mathbf{w}, b, \xi} C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s.t. } \max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) = \xi_n, \quad n \in \{1, \dots, N\}$$

Recovering our previous SVM formulation

Minimizing hinge loss

$$\min_{\mathbf{w}, b, \xi} C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2$$

s.t. $\max(0, 1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b]) = \xi_n, \quad n \in \{1, \dots, N\}$

$(\mathbf{w}^*, b^*), \{\xi_n^*\}$

is equivalent to

SVM

$$\min_{\mathbf{w}, b, \xi} \left[C \sum_n \xi_n + \frac{1}{2} \|\mathbf{w}\|_2^2 \right]$$

s.t. $1 - y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \leq \xi_n, \quad n \in \{1, \dots, N\}$

$0 \leq \xi_n, \quad n \in \{1, \dots, N\}$

$\xi_n^* = \max(0, 1 - y_n(\mathbf{w}^{*T} \phi(\mathbf{x}_n) + b^*))$

$\Rightarrow \xi_n \geq \max(0, 1 - y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b))$

$\Rightarrow \xi_n^* \geq \max(0, 1 - y_n(\mathbf{w}^{*T} \phi(\mathbf{x}_n) + b^*))$

At optimal solution constraints are active so we have equality! Why?

$$\underline{\underline{\xi_n^*}} > \max \left(0, 1 - y_n (\omega^{*T} \phi(x_n) + b^*) \right)$$

$$\underline{\underline{\xi_n^2}}$$

$$C \sum_n \xi_n + \frac{1}{2} \|\omega\|_2^2$$

↑

At optimal solution constraints are active so we have equality! Why?

- If $\xi_n^* > \max(0, 1 - y_n f(\mathbf{x}_n))$, we could choose $\bar{\xi}_n < \xi_n^*$ and still satisfy the constraint while reducing our objective function!
- Since $c \geq \max(a, b) \iff c \geq a, c \geq b$, we recover previous formulation

SVM: hyperplane with largest margin
= minimizes hinge loss
+ ℓ_2 -regularization

Outline

- 1 Review of last lecture
- 2 SVM – Hinge loss
- 3 Kernelized SVMs**
- 4 Evaluating ML algorithms

Goal: Kernelize SVMs

$$\underline{\phi(\mathbf{x}_n)}$$

- Rewrite the SVM optimization problem so that it no longer depends on $\phi(\mathbf{x}_n)$ but instead depends only on inner products $\phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$.
- If we can do this, we can then replace $\phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$ with a kernel function $k(\mathbf{x}_n, \mathbf{x}_m)$.

Primal formulation of support vector machines (SVM)

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \quad \text{and} \quad \xi_n \geq 0, \quad n \in \{1, \dots, N\} \end{aligned}$$

This can be converted to its dual form.

Dual formulation of SVM

Let w_* be the minimizer of the SVM problem for training data: $\{(\mathbf{x}_n, y_n)\}$.

Then w_* = $\sum_n \alpha_n y_n \phi(\mathbf{x}_n)$ for a new set of dual variables: $\alpha_n \geq 0$.

Primal (arrow pointing to w_*)

Dual (arrow pointing to α_n)

Dual formulation of SVM

$$y_m \in \{+1, -1\}$$

Dual is also a convex program

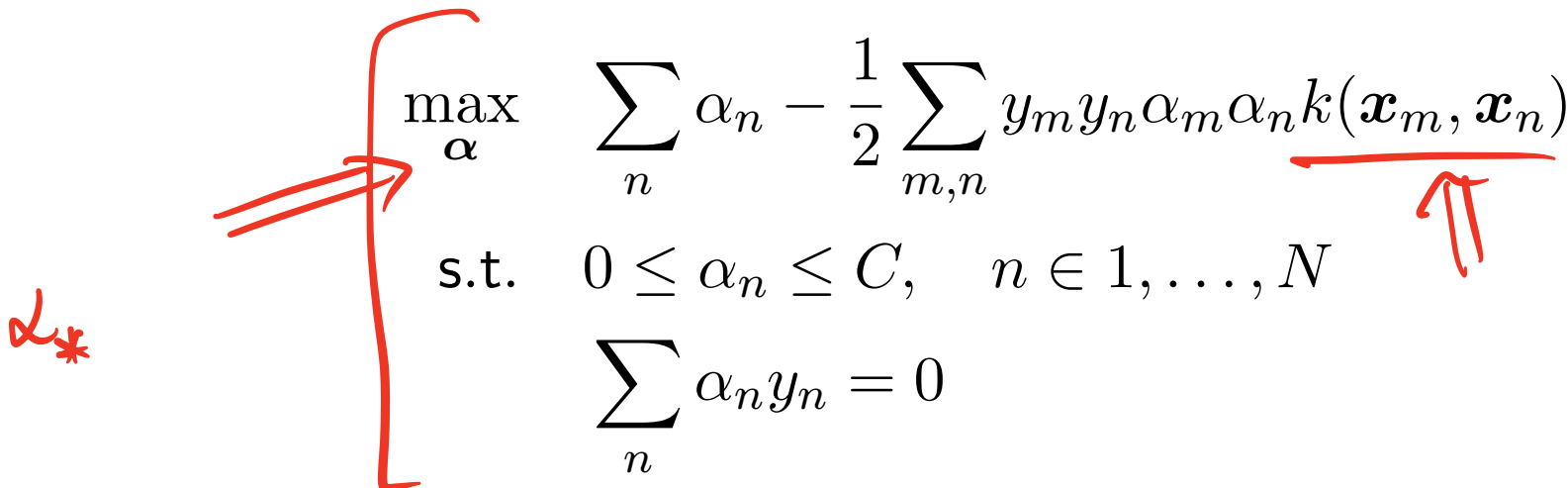
$$\begin{aligned} \max_{\alpha} & \left(\sum_n \alpha_n - \frac{1}{2} \sum_{m,n} \underline{y_m y_n} \alpha_m \alpha_n \underline{\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)} \right) \\ \text{s.t.} & \quad 0 \leq \alpha_n \leq \underline{C}, \quad n \in 1, \dots, N \\ & \quad \sum_n \alpha_n y_n = 0 \end{aligned}$$

$\alpha = (\alpha_1, \dots, \alpha_N)$

There are N dual variable α_n , one for each constraint in the primal formulation

Kernel SVM

We replace the inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ with a kernel function


$$\begin{aligned} \max_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \underline{k(\mathbf{x}_m, \mathbf{x}_n)} \\ \text{s.t.} \quad & 0 \leq \alpha_n \leq C, \quad n \in 1, \dots, N \\ & \sum_n \alpha_n y_n = 0 \end{aligned}$$

We can define a kernel function to work with nonlinear features and learn a nonlinear decision surface

Recovering solution to the primal formulation

$$\underline{\underline{w}}^T \underline{\underline{\phi(x_n)}} + \underline{\underline{b}}$$

$$\underline{\underline{k(x, x_n)}}$$

Weights w (primal variables) relation to α (dual variables)

$$\underline{\underline{w}} = \sum_n y_n \underline{\underline{\alpha_n}} \underline{\underline{\phi(x_n)}} \leftarrow \text{Linear combination of the input features}$$

Prediction on a test point x

$$h(\underline{\underline{x}}) = \text{SIGN}(\underline{\underline{w}}^T \underline{\underline{\phi(x)}} + \underline{\underline{b}}) = \text{SIGN}(\underbrace{\sum_n y_n \alpha_n k(\underline{\underline{x_n}}, \underline{\underline{x}})}_{\text{kernel function}} + \underline{\underline{b}})$$

At test time it suffices to know the kernel function!

$$\underline{\underline{w}} = \underline{\underline{\sum_n y_n \alpha_n \phi(x_n)}}$$

$$\begin{aligned} \underline{\underline{w}}^T \underline{\underline{\phi(x)}} &= \left(\sum_n y_n \underline{\underline{\alpha_n}} \underline{\underline{\phi(x_n)}} \right)^T \underline{\underline{\phi(x)}} \\ &= \sum_n y_n \alpha_n \underline{\underline{\phi(x_n)^T \phi(x)}} \end{aligned}$$

Recovering solution to the primal formulation

We already identified the primal variable \mathbf{w} as

$$\underline{\mathbf{w}} = \sum_n \alpha_n y_n \underline{\phi(\mathbf{x}_n)}$$

- From the dual, we know that $0 \leq \alpha_n \leq C$.

$$0 \leq \alpha_n \leq C$$

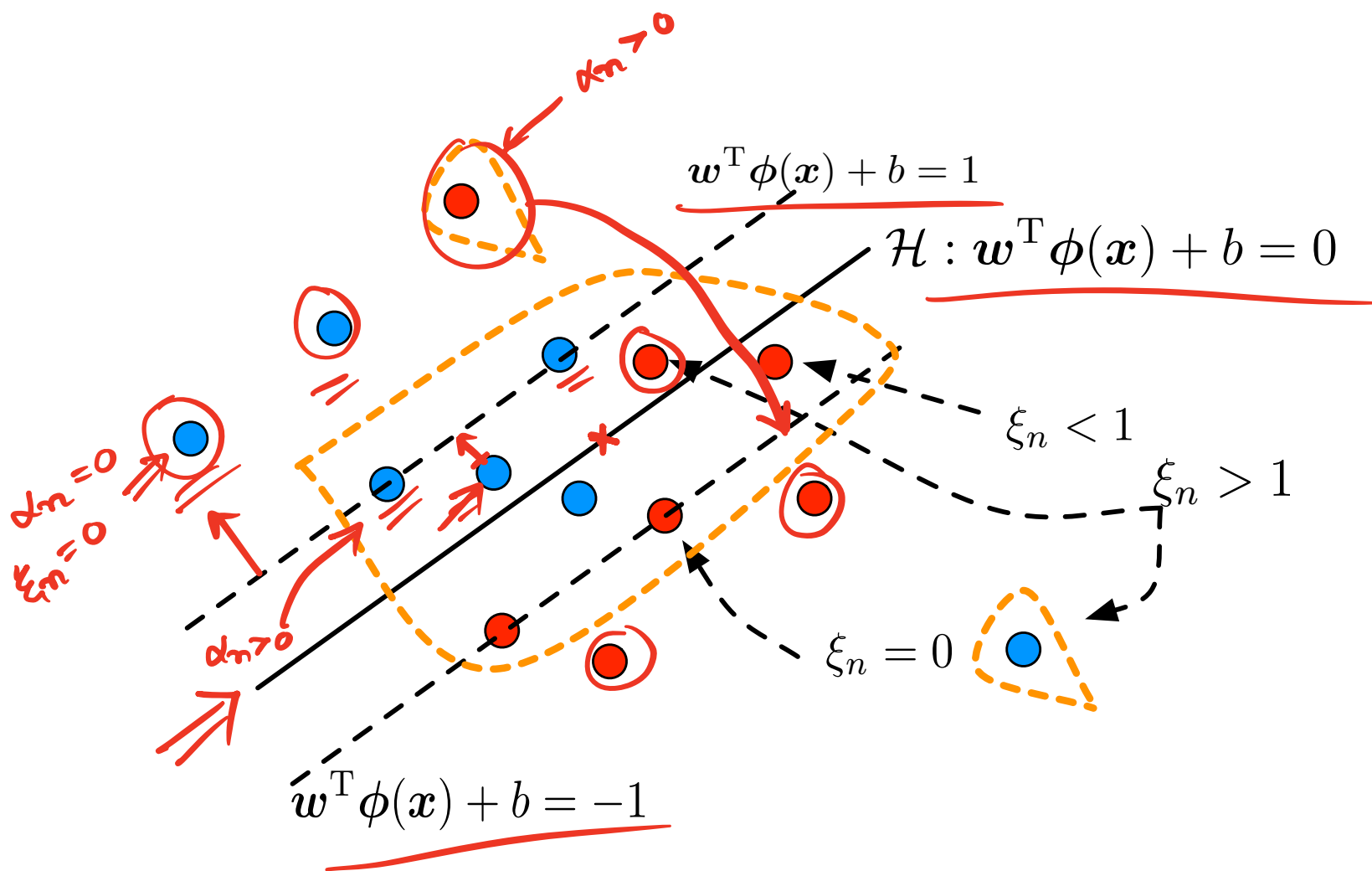
- If $\alpha_n = 0$, then the example n does not affect the hyperplane/decision boundary computed by SVM.

$$y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1.$$

- Only those examples with $\alpha_n > 0$ affect the hyperplane/decision boundary. These examples are termed **support vectors**.

- When will $\alpha_n > 0$?
 $\xi_n \geq 0$
 - $\alpha_n = C \Rightarrow y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] \leq 1.$
 - $0 < \alpha_n < C \Rightarrow y_n[\mathbf{w}^T \phi(\mathbf{x}_n) + b] = 1.$

Visualization of how training data points are categorized



Support vectors are highlighted by the dotted orange lines

Outline

- 1 Review of last lecture
- 2 SVM – Hinge loss
- 3 Kernelized SVMs
- 4 Evaluating ML algorithms**

Binary classification

Goal is to achieve high accuracy

Accuracy: defined to be the mean of the 0/1 loss $l(y, \hat{y})$.

$$\ell(y, h(\mathbf{x})) = \mathbf{1}\{y \neq h(\mathbf{x})\}$$
$$\underline{\underline{R^{\text{EMP}}[\underline{\underline{h(\mathbf{x})}]} = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)}}$$

Assumes $y \in \{0, 1\}$ and h outputs values in $\{0, 1\}$.

Binary classification

Goal is to achieve high accuracy

Accuracy: defined to be the mean of the 0/1 loss $\ell(y, \hat{y})$.

$$\ell(y, h(\mathbf{x})) = \mathbf{1}\{y \neq h(\mathbf{x})\}$$

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Assumes $y \in \{0, 1\}$ and h outputs values in $\{0, 1\}$.

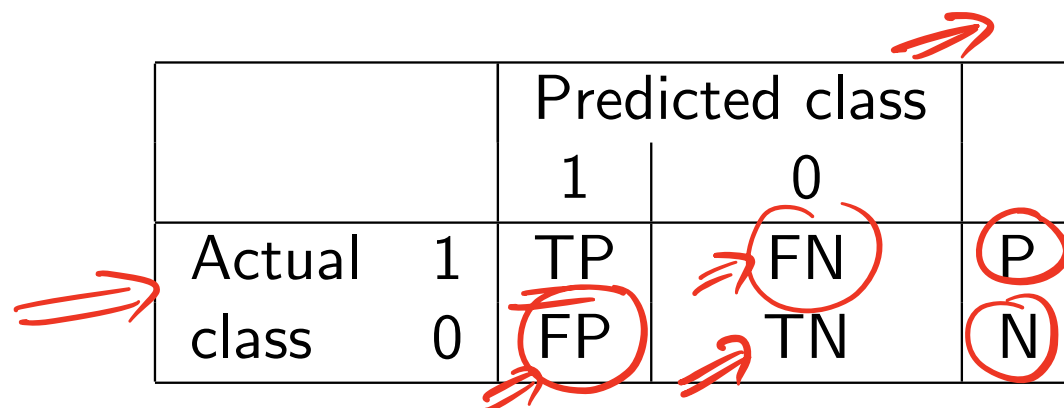
May not be what we want

Is this email spam?

Does the patient have cancer?

Confusion matrix

Given a dataset of P positive instances and N negative instances,



The diagram shows a confusion matrix with red annotations. Red arrows point to the 'Actual class' label, the 'Predicted class' header, and the '1' and '0' headers. Red circles highlight the 'FN' and 'FP' cells, and the 'P' and 'N' labels in the rightmost column.

		Predicted class		
		1	0	
Actual class	1	TP	FN	P
	0	FP	TN	N

P = Positive, N = Negative

TN = True negative, FP = False positive

TP = True positive, FN = False negative

Spam classification

- Positive : spam, Negative: non-spam
- TP: Spam email classified as spam
- FP: Non-spam email classified as spam

Confusion matrix

Given a dataset of P positive instances and N negative instances,

		Predicted class		
		1	0	
Actual class	1	TP	FN	P
	0	FP	TN	N

P = Positive, N = Negative

TN = True negative, FP = False positive

TP = True positive, FN = False negative

- Accuracy = $\frac{TP + TN}{P + N}$

- Error = $1 - \text{Accuracy}$

Confusion matrix

Given a dataset of P positive instances and N negative instances,

		Predicted class		
		1	0	
Actual class	1	TP	FN	P
	0	FP	TN	N

P = Positive, N = Negative

TN = True negative, FP = False positive

TP = True positive, FN = False negative

- True Positive Rate (TPR) = Recall = Sensitivity = $\frac{TP}{P}$
- Probability of classifying a spam email as spam.

Confusion matrix

Given a dataset of P positive instances and N negative instances,

		Predicted class		
		1	0	
Actual class	1	TP	FN	P
	0	FP	TN	N

P = Positive, N = Negative

TN = True negative, FP = False positive

TP = True positive, FN = False negative

- Specificity = $\frac{TN}{N}$
- Probability of classifying a non-spam email as non-spam.
- False Positive Rate (FPR) = 1 - Specificity

Confusion matrix

Given a dataset of P positive instances and N negative instances,

		Predicted class		
		1	0	
Actual class	1	TP	FN	P
	0	FP	TN	N

P = Positive, N = Negative

TN = True negative, FP = False positive

TP = True positive, FN = False negative

- Precision = Positive Predictive Value (PPV) = $\frac{TP}{TP+FP}$
- Probability that an email classified as spam is spam.
- Different from specificity!

Visualizing performance

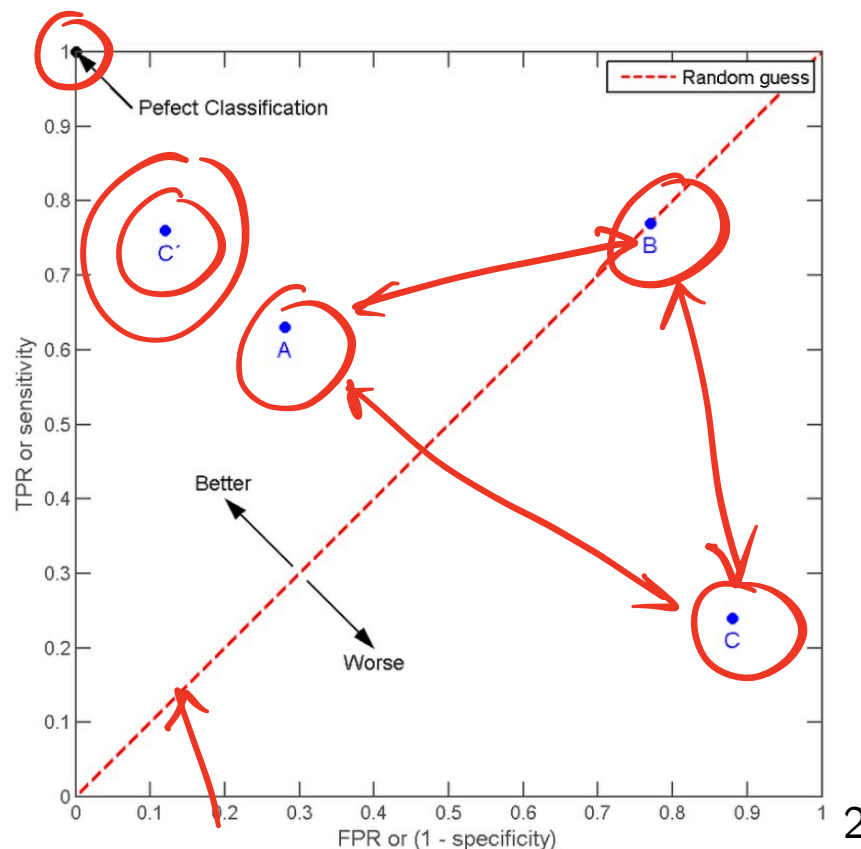
Many different performance measures

A			B			C			C'		
TP=63	FN=37	100	TP=77	FN=23	100	TP=24	FN=76	100	TP=76	FN=24	100
FP=28	TN=72	100	FP=77	TN=23	100	FP=88	TN=12	100	FP=12	TN=88	100
91	109	200	154	46	200	112	88	200	88	112	200

1

ROC curve

TPR (sensitivity) vs FPR (1-specificity)



- Best possible method is in the upper left corner.
- Random guess would fall on the diagonal.

Sometimes we care about the confidence of predictions

- For linear models (e.g., perceptron, logistic regression), the activation $a(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ gives us a ranking of the predictions.
- In the case of the perceptron, given $\underline{x_1}$ and $\underline{x_2}$ such that $\underline{a(x_1)} > \underline{a(x_2)} > 0$.
 - ▶ Both $\underline{x_1}$ and $\underline{x_2}$ are classified as positive.
 - ▶ Perceptron is more confident in the prediction for x_1 than for x_2 .

ROC curve allows us to visualize the ranking of predictions

Building a ROC curve

- Sort the instances according to the confidence of instance being in class 1 (positive class).
 - ▶ How confident is the classifier that an email is spam?

Building a ROC curve

- Sort the instances according to the confidence of instance being in class 1 (positive class).
 - ▶ How confident is the classifier that an email is spam?
- Step through the sorted list from high to low confidence.
- As you step through, build a binary classifier that predicts the top k instances as belonging to class 1 and the remaining to class 0.
 - ▶ How aggressively do you want to mark a message as spam ?

Building a ROC curve

- Sort the instances according to the confidence of instance being in class 1 (positive class).
 - ▶ How confident is the classifier that an email is spam?
- Step through the sorted list from high to low confidence.
- As you step through, build a binary classifier that predicts the top k instances as belonging to class 1 and the remaining to class 0.
 - ▶ How aggressively do you want to mark a message as spam ?
- Compute the TPR (True Positive Rate) and FPR (False Positive Rate).
 - ▶ High up in the list: Low TPR and low FPR.
 - ▶ As we move down, TPR increases but so does FPR.

Building a ROC curve

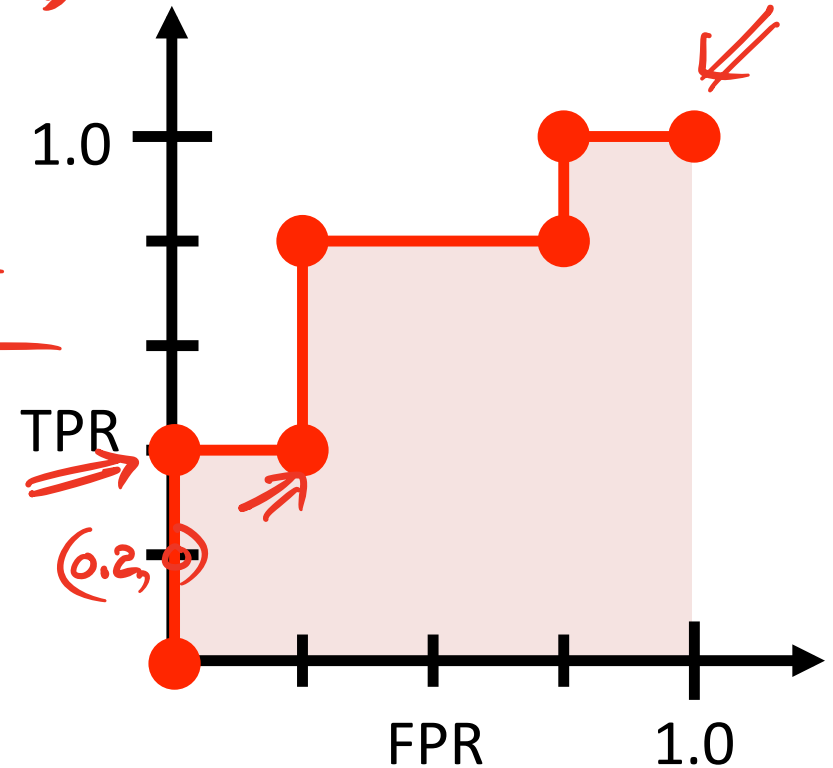
- Sort the instances according to the confidence of instance being in class 1 (positive class).
 - ▶ How confident is the classifier that an email is spam?
- Step through the sorted list from high to low confidence.
- As you step through, build a binary classifier that predicts the top k instances as belonging to class 1 and the remaining to class 0.
 - ▶ How aggressively do you want to mark a message as spam ?
- Compute the TPR (True Positive Rate) and FPR (False Positive Rate).
 - ▶ High up in the list: Low TPR and low FPR.
 - ▶ As we move down, TPR increases but so does FPR.
- Plot TPR vs FPR.

Example

$$TPR = \frac{TP}{Positives} = \frac{1}{5}$$

$$FPR = 1 - \frac{TN}{Negatives} = 1 - \frac{5}{5} = 0$$

Instance	Confidence in positive	Actual class
9	0.99	1
7	0.98	1
1	0.72	0
2	0.70	1
6	0.65	1
10	0.51	0
3	0.39	0
5	0.24	1
4	0.11	0
8	0.01	0



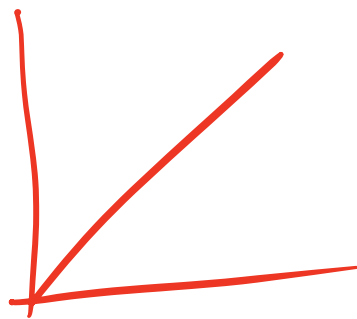
$$TPR = \frac{2}{5}$$

$$FPR = \frac{1}{5}$$

Interpreting ROC curves

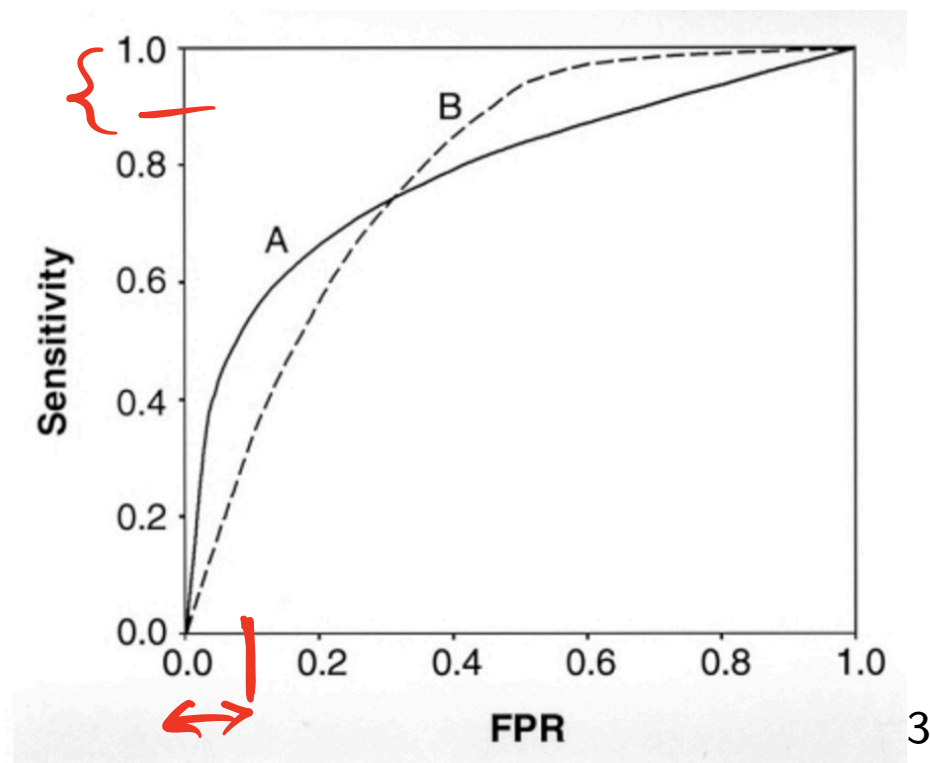
Area under the ROC Curve (AUC or AUROC)

- The probability that a classifier will rank a randomly chosen positive instance (instance with class label 1) higher than a randomly chosen negative one (instance with class label 0).
- A measure of the overall performance of a classifier.
- A random classifier has $\text{AUROC} = \frac{1}{2}$.
- A perfect classifier has $\text{AUROC} = 1$.



Interpreting ROC curves

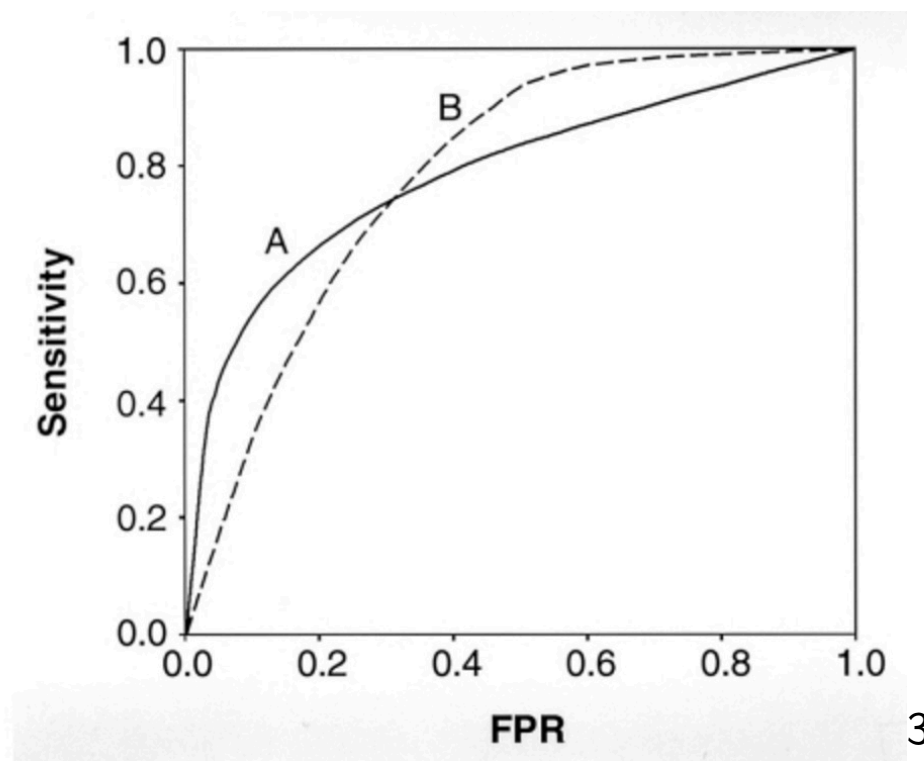
Classifiers A and B with equal AUROC



Which one do we choose?

Interpreting ROC curves

Classifiers A and B with equal AUROC



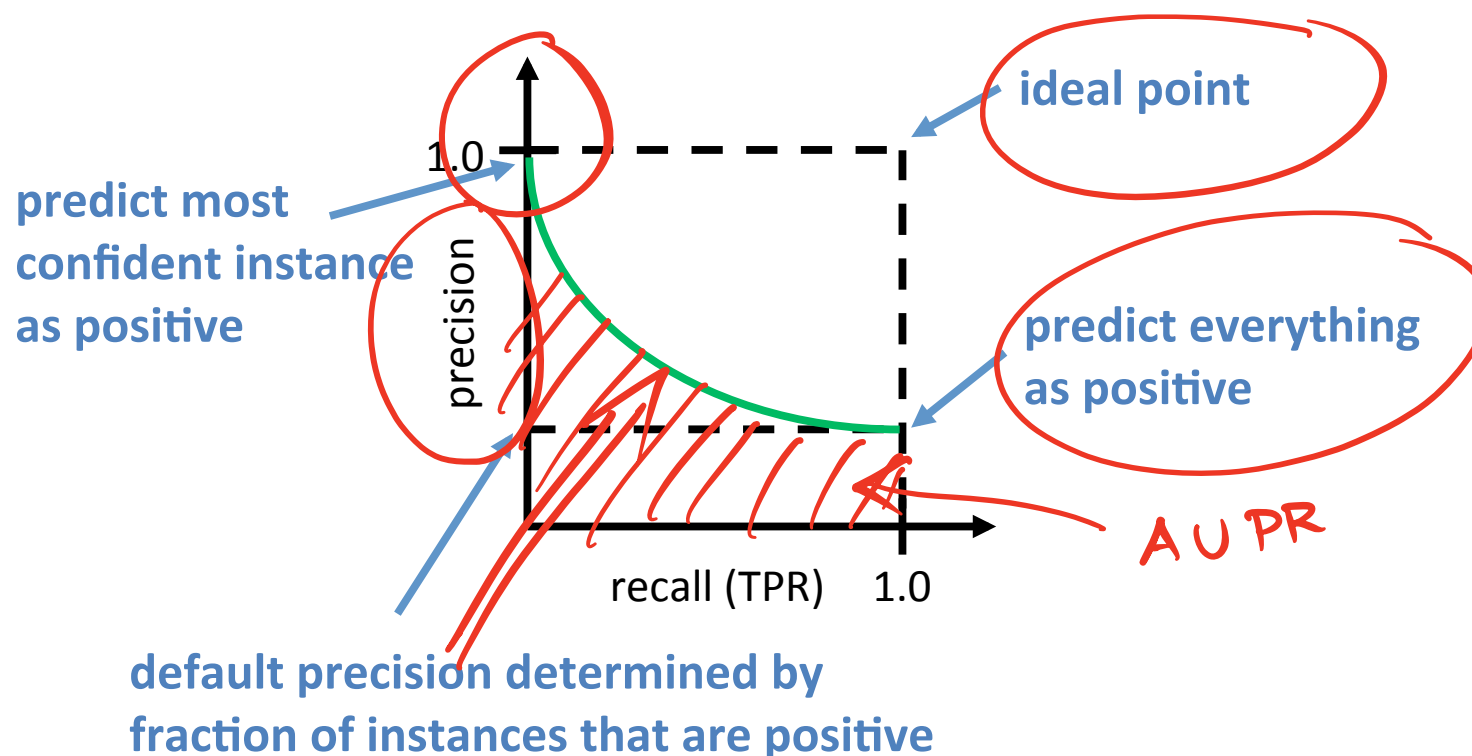
Which one do we choose?

Depends on whether we care about sensitivity or FPR

Other trade-offs

$$\text{Precision} = \frac{TP}{TP + FP}$$

Precision-recall curve



Interpreting precision-recall curves

Area under the precision-recall curve (AUPR)

- Can compute the area under the precision-recall curve.

Interpreting precision-recall curves

F₁-score

$$F_1 = \frac{1}{\frac{1}{2} \frac{1}{\text{PRECISION}} + \frac{1}{2} \frac{1}{\text{RECALL}}}$$

- Want F_1 to be high.
- The Harmonic mean of precision and recall.
- A single number that trades off precision and recall.
- If a classifier has high precision and low recall (or other way round), its F_1 -score will be low.
- Encourages precision and recall values that are similar.

Summary

- Many ways to measure accuracy.
- Think carefully about which ones are most relevant.
- Often there are trade-offs
 - ▶ TPR vs FPR, Precision vs recall
- Graphical displays useful.
 - ▶ ROC and Precision-recall curves display performance at various levels of confidence.
 - ▶ Can be summarized by area under the curve (AUROC, AUPR) as well as numbers such as the F_1 -score.