CSCE 978 Lecture 3: Risk and Loss Functions*	Introduction
Stephen D. Scott	 In Lecture 1 we mentioned our desire to infer a "good" classifier What does this mean?!?!
January 24, 2006	• There are many ways to define "goodness", even for binary classification
*Most figures ©2002 MIT Press, Bernhard Schölkopf, and Alex Smola. 1	2
	Loss Functions
Outline Loss functions Binary classification Regression 	D3.1 Let $(x, y, f(x)) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$ be the pattern x , its true label y and a prediction $f(x)$ of y . A loss function is a mapping $c : \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ with the property $c(x, y, y) = 0$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ • c is always ≥ 0 so we can't use good predictions to "undo" bad ones
 Expected risk Sections 1.3, 3.1–3.2 (also read Section 3.5) 	 It is always possible to get 0 loss on pattern x by predicting correctly Our choice of loss function will depend on con-

Loss Functions

Binary Classification

• Count number of misclassifications:

$$c(x, y, f(x)) = \begin{cases} 0 & \text{if } y = f(x) \\ 1 & \text{otherwise} \end{cases}$$

• Same as above, but penalty is input-dependent:

$$c(x, y, f(x)) = \begin{cases} 0 & \text{if } y = f(x) \\ \tilde{c}(x) & \text{otherwise} \end{cases}$$

- E.g. if $y \in \{\text{rocks}, \text{diamonds}\}$ then penalty for "false diamond" classification depends on x's weight
- Can also have different values for false positive (y = -1) and false negative (y = +1) errors
 - If $y \in \{\text{cancer}, \neg \text{cancer}\}$ then FP results in unnecessary treatment, but FN can be fatal

Loss Functions Binary Classification (cont'd)

- If f(x) is real-valued and $y \in \{-1, +1\}$, can think of sign(f(x)) as prediction and |f(x)| as a confidence. Then a highly confident incorrect prediction can be penalized more, as can low-confidence correct predictions:
 - Soft margin loss:

$$c(x, y, f(x)) = \max(0, 1 - yf(x))$$
$$= \begin{cases} 0 & \text{if } yf(x) \ge 1\\ 1 - yf(x) & \text{otherwise} \end{cases}$$

- Logistic loss:

$$c(x, y, f(x)) = \ln \left(1 + \exp(-yf(x))\right)$$

 Both penalize a lot for confident, incorrect predictions, penalize a little for low confidence, and don't penalize much or at all for confident, correct predictions

6

Loss Functions Regression

- In regression, $\mathcal{Y}\subseteq\mathbb{R}$ rather than $\mathcal{Y}=\{-1,+1\}$
- Thus we're interested in how far off our prediction f(x) is
- Squared loss (very popular):

$$c(x, y, f(x)) = (f(x) - y)^2$$

• Can extend soft margin loss to ϵ -insensitive loss, which doesn't penalize for close predictions:

$$c(x, y, f(x)) = |f(x) - y|_{\epsilon} = \max(|f(x) - y| - \epsilon, 0)$$





0

y f(x)

0

y f(x)

0

y f(x)

Loss Functions

7

y f(x)

5

8

	Risk
Loss Functions Practical Considerations	 A loss function measures error on individual examples Our ultimate goal is to minimize loss on new (yet unseen) examples
ant loss function to be:	How do we measure this?
Cheap to compute	 Without making certain assumptions, this is very difficult or even impossible
Have few discontinuities in first derivative	- Assume that there is a probability distribu- tion $P(x, y)$ on $\mathcal{X} \times \mathcal{Y}$ that governs genera-
Convex (to ensure unique global optimum)	tion of patterns and labels
Yield computationally efficient solutions for learning	* Assume the pairs (x, y) are drawn iid (in- dependent and identically distributed) ac- cording to $P(x, y)$
Resistant to outliers/noise	* Generally, we won't make specific assumptions about the nature of $P(x,y)$
	- $P(y x) =$ conditional probability of getting label y given that x is the pattern (so x could have a different label on each draw)
9	10
Risk Definitions	Risk Definitions (cont'd)

• For now, assume we know all the new patterns we'll ever classify; call these the test patterns $x_1',\ldots,x_{m'}'$ (note we do not know the labels until after we make predictions)

Ŵ

D3.2 When test set $x'_1, \ldots, x'_{m'}$ already known, goal is to minimize the <u>expected error on the test set</u>:

$$R_{\mathsf{test}}[f] := \frac{1}{m'} \sum_{i=1}^{m'} \int_{\mathcal{Y}} c(x'_i, y, f(x'_i)) \, d\mathsf{P}(y \mid x'_i)$$

- Often, minimizing $R_{\text{test}}[f]$ not realistic since typically don't know test set a priori
 - One exception: querying fixed collection of images, biological sequences, etc.

D3.3 The expected risk (expected loss) wrt P & c:

$$R[f] := \mathbf{E} [R_{\text{test}}[f]] = \mathbf{E} [c(x, y, f(x))]$$
$$= \int_{\mathcal{X} \times \mathcal{Y}} c(x, y, f(x)) \, d\mathbf{P}(x, y)$$

• Not realistic since we don't know P(x, y)

- To get a handle on P(x, y), assume it's the same one that generated the training set
- Now use the training patterns to estimate P(x, y)

D3.4 The empirical risk is

$$R_{emp}[f] := \int_{\mathcal{X} \times Y} c(x, y, f(x)) p_{emp}(x, y) \, dx \, dy$$
$$= \frac{1}{m} \sum_{i=1}^{m} c(x_i, y_i, f(x_i))$$

- Easy to compute and generally straightforward to minimize (depending on c)
- So now all we have to do is find an f that minimizes $R_{emp}[f]$, use that as our predictor, and we're done, right?

(Can we go home now?)

12

 NO! We have to appropriately <u>restrict</u> the set of functions <i>F</i> from which we choose <i>f</i> Otherwise, <i>R</i>_{emp}[<i>f</i>] won't approximate <i>R</i>[<i>f</i>], which is what we want to minimize E.g. what if <i>F</i> is the set of all functions from <i>X</i> to <i>Y</i>? Then our learning algorithm could get <i>R</i>_{emp}[<i>f</i>] = 0 by simply storing the (<i>x</i>, <i>y</i>) pairs in a table (i.e. memorization) Is this learning? Will it generalize well? 	Topic summary (over Lectures 2 and 3) due in 1 week!
 Restricting <i>F</i> has been looked from many perspectives: e.g. VC dimension, bias, structural risk minimization Our approach (called <u>regularization</u>) will quantify the "power" ("expressiveness") of each <i>f</i> and minimize a sum of this and <i>R</i>_{emp}[<i>f</i>] Special case: <u>minimum description length</u> principle 	14