

Classification

Sewoong Oh

CSE/STAT 416

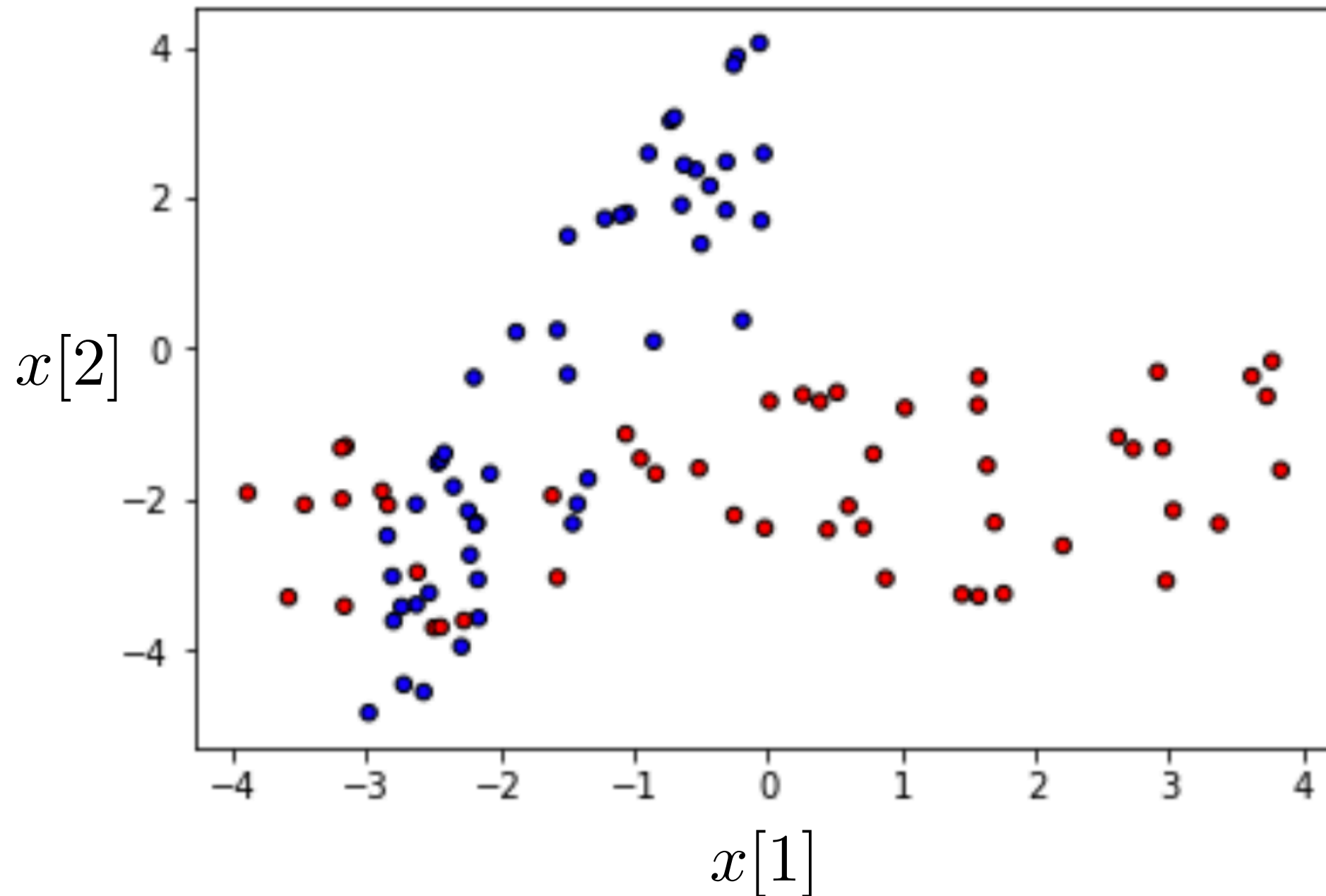
University of Washington

Boolean Classification

Boolean classification

- **Supervised learning** is training a predictor from labelled examples:
- There are two types of supervised learning
 - 1. **Regression**: the output variable \mathbf{y} to be predicted is **real valued** scalar or a vector
 - 2. **Classification**: the output variable \mathbf{y} to be predicted is categorical
 - 2.1 **Boolean classification**: there are two classes
 - 2.2 **Multi-class classification**: multiple classes
- We study Boolean classification in this chapter
- We denote two classes by -1 and 1, often corresponding to {FALSE, TRUE}
- for a data point (x_i, y_i) , the value $y_i \in \{-1, 1\}$ is called the **class** or **label**
- A **Boolean classifier** predicts label \mathbf{y} given input \mathbf{x}

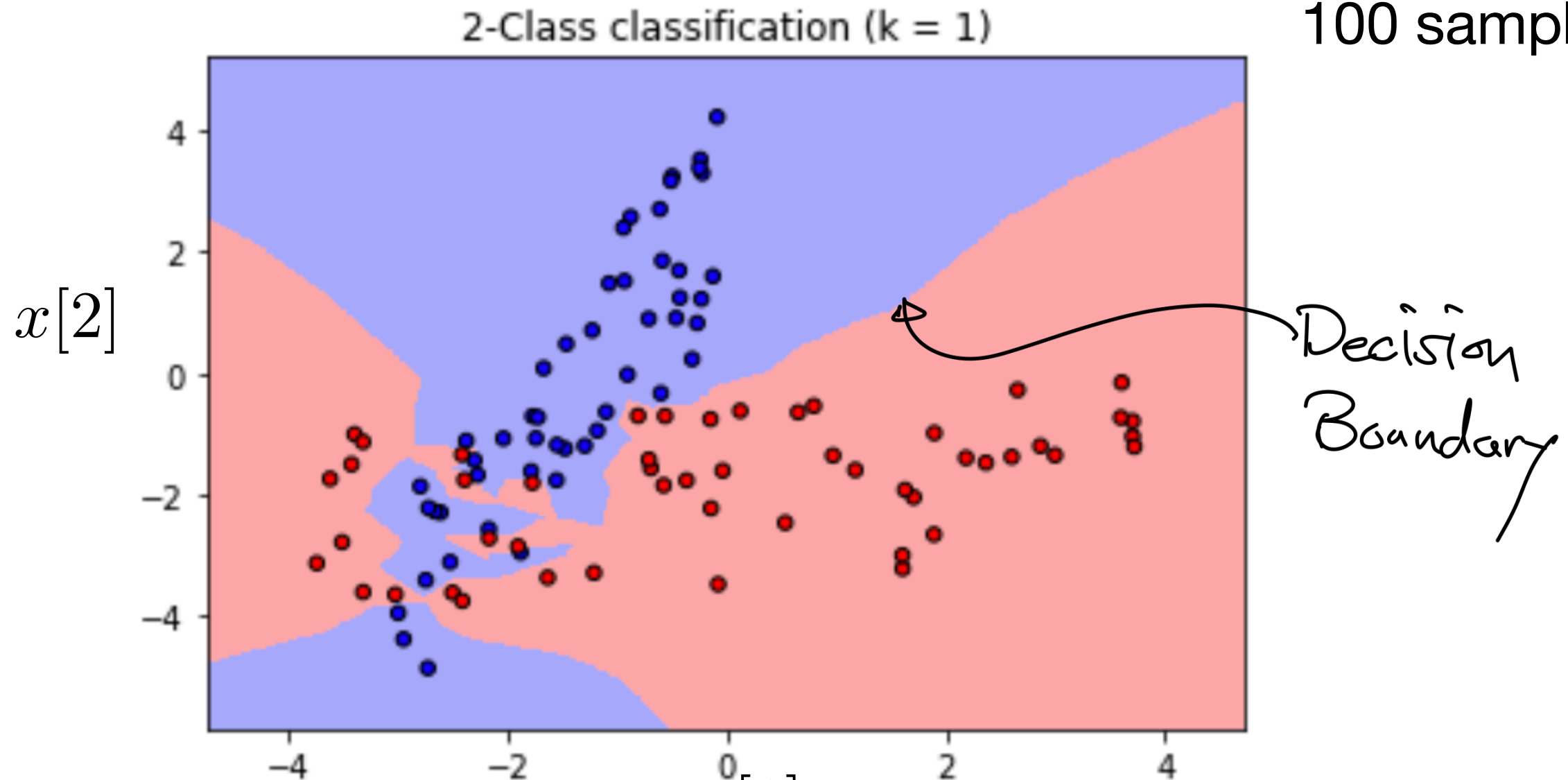
Classification



- in this example $x_i \in \mathbf{R}^2$
- Red points have label $y_i = -1$, blue points have label $y_i = 1$
- We want a predictor that maps any \mathbf{x} into prediction $\hat{y} \in \{-1, 1\}$

Example: nearest neighbor

- Trained on 100 samples



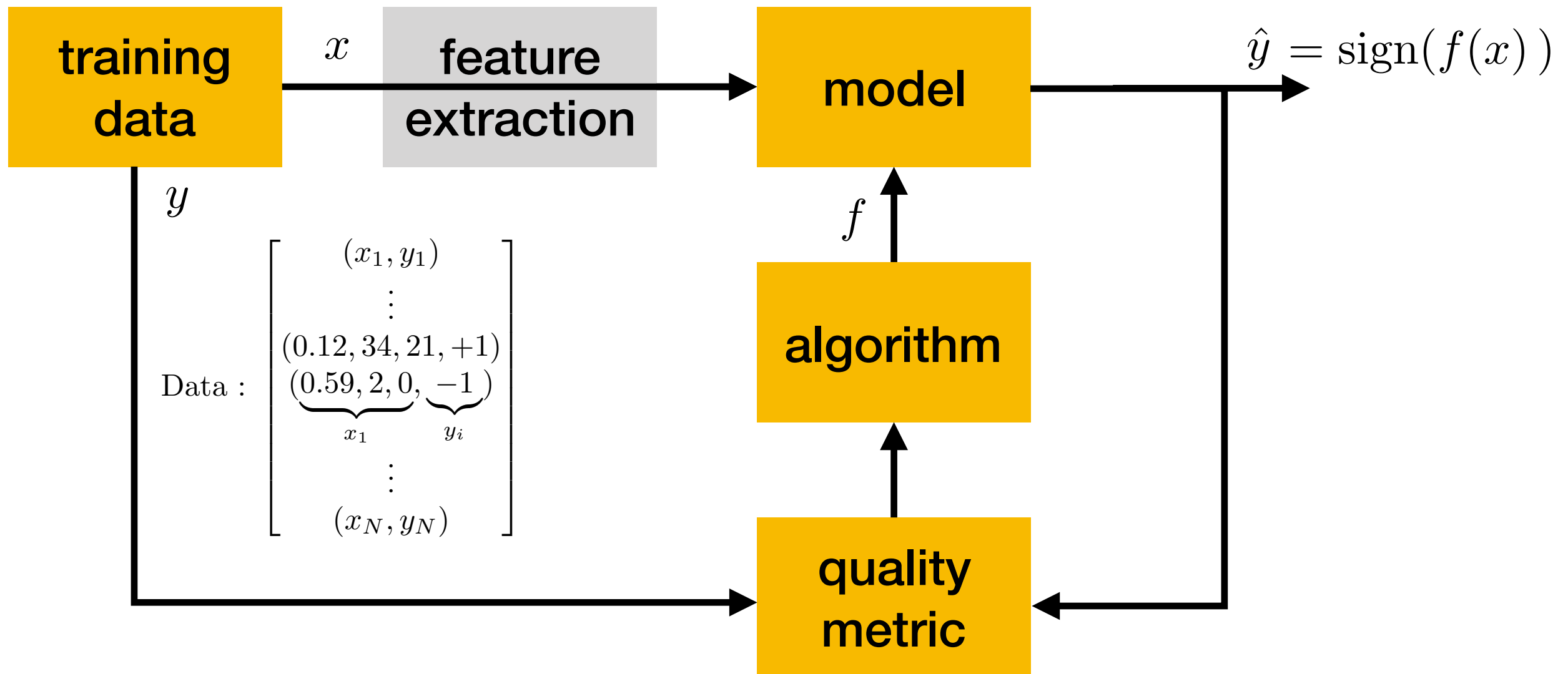
- given x , let $k = \arg \min_i \|x - x_i\|$, then predict $\hat{y} = y_k$
- Red region is the set of x for which prediction is -1
- Blue region is the set of x for which prediction is 1
- Zero training error, but overfitting

Simple approach: Linear classifier

linear model: $f(x) = w_0 + w_1x[1] + w_2x[2] + \dots$

In training: find w that minimizes $\frac{1}{N} \sum_{i=1}^N (f(x_i) - y_i)^2$

In prediction: $\hat{y} = \text{sign}(f(x_i))$



- This is the simplest form: L2 regression followed by $\text{sign}()$
- General recipe: (some type of) regression followed by $\text{sign}()$
- We can do significantly better by choosing the right loss

Example: sentiment analysis

List of positive words	List of negative words
great, awesome, good, amazing,...	bad, terrible, disgusting, sucks,...

x_i



Sushi was great, the food was awesome, but the service was terrible.



$$\hat{y}_i = \text{sign}(\text{number of positive words} - \text{number of negative words})$$

- If we have access to the list of positive and negative words, then we could count them to give a score $f(x)$ and take the sign for estimating the sentiment in {positive,negative}

Example: sentiment analysis

- **Linear classifier**

Sushi was great, the food was awesome, but the service was terrible.

$h_j(x)$ = how many times the word appears
 w_j = how positive is that word

→ $\hat{y}_i = \text{sign}(w_0 + w_1 h_1(x) + w_2 h_2(x) + \dots)$

Word	Weight
good	1.0
great	1.5
awesome	2.7
bad	-1.0
terrible	-2.1
awful	-3.3
restaurant, the, we, where, ...	0.0
...	...

- Without manually constructed list, we can use ML to learn the sentiment of the words (parameters \mathbf{w}), and then compute a score $f(x)$

Confusion Matrix:
**How should we measure performance
in Boolean classification?**

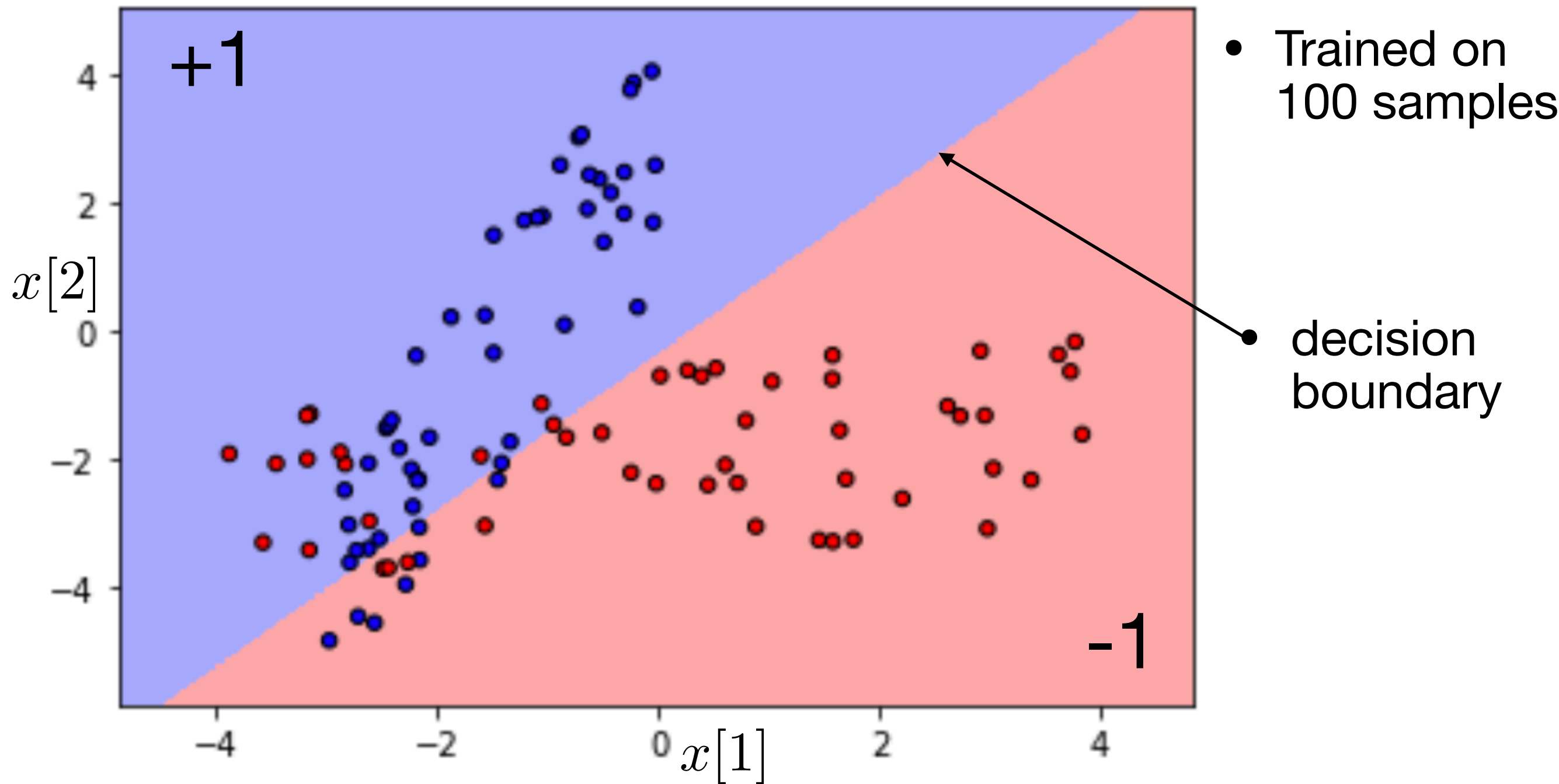
Two types of error

- When measuring performance of a predictor on Boolean classification
- each input data x_i has a label $y_i \in \{-1, 1\}$
- each corresponding prediction is $\hat{y}_i \in \{-1, 1\}$
- Only four possible combinations of (\hat{y}_i, y_i) :
 - **true positive** if $\hat{y}_i = 1$ and $y_i = 1$
 - **true negative** if $\hat{y}_i = -1$ and $y_i = -1$
 - **false negative of type II error** if $\hat{y}_i = -1$ and $y_i = 1$
 - **false positive of type I error** if $\hat{y}_i = 1$ and $y_i = -1$

$$C = \begin{bmatrix} \# \text{ true negatives} & \# \text{ false negatives} \\ \# \text{ false negatives} & \# \text{ true positives} \end{bmatrix} = \begin{bmatrix} C_{tn} & C_{fn} \\ C_{fp} & C_{tp} \end{bmatrix}$$

- Each outcome of a predictor means different things depending on the application: Screening for cancer...

Example: linear classifier



- 18% mis-classified in training data
- true positive $C_{tp} = 42$, false positive $C_{fp} = 12$,
- true negative $C_{tn} = 38$, false negative $C_{fn} = 8$

We can represent the performance with a confusion matrix

- Define the **confusion matrix**: (some people use the transpose of C)

$$\begin{bmatrix} C_{tn} & C_{fn} \\ C_{fp} & C_{tp} \end{bmatrix} = \begin{bmatrix} 38 & 8 \\ 12 & 42 \end{bmatrix}$$

- $C_{tn} + C_{fn} + C_{fp} + C_{tp} = N$
- $N_n = C_{tn} + C_{fp}$ is the number of negative examples
- $N_p = C_{fn} + C_{tp}$ is the number of positive examples
- Diagonal entries give numbers of correct predictions
- Off-diagonal entries give numbers of incorrect predictions

Some Boolean classification measures

- **Confusion matrix**
$$\begin{bmatrix} C_{tn} & C_{fn} \\ C_{fp} & C_{tp} \end{bmatrix} = \begin{bmatrix} 38 & 8 \\ 12 & 42 \end{bmatrix}$$

- The basic error measures are:

- **False positive rate** is C_{fp}/N (e.g. False discovery rate)

- **False negative rate** is C_{fn}/N (e.g. medical diagnosis)

- **Error rate** is $(C_{fn} + C_{fp})/N$

- **Accuracy** is $(C_{tn} + C_{tp})/N$

- High accuracy does not always mean good classifier

- For example, 99% population does not have cancer, and predicting always no cancer achieves accuracy 99%

$$C = \begin{bmatrix} 99 & 1 \\ 0 & 0 \end{bmatrix}$$

- Error measures also used:

- **True positive rate** or **sensitivity** or **recall** is C_{tp}/N_p (e.g. 84%, 0%)

- **False alarm rate** is C_{fp}/N_n (e.g. 24%, 0%)

- **Specificity** or **true negative rate** is C_{tn}/N_n (e.g. 76%, 100%)

- **Precision** is $C_{tp}/(C_{tp} + C_{fp})$ (e.g. 77%, 0%)

Neyman-Pearson error

- **Neyman-Pearson error** over a data set is

$$\kappa C_{fn}/N + C_{fp}/N$$

- A scalarization of our two objectives, minimizing false positive and minimizing false negative rates
- A positive real values k is how much more false negative irritates us than false positives
- When $k=1$, the Neyman-Pearson error is the **error rate**
- A common and flexible measure of error
- Baseline for designing a loss function for training
 - Similar to how MSE performance metric is a baseline for designing loss for regression problems

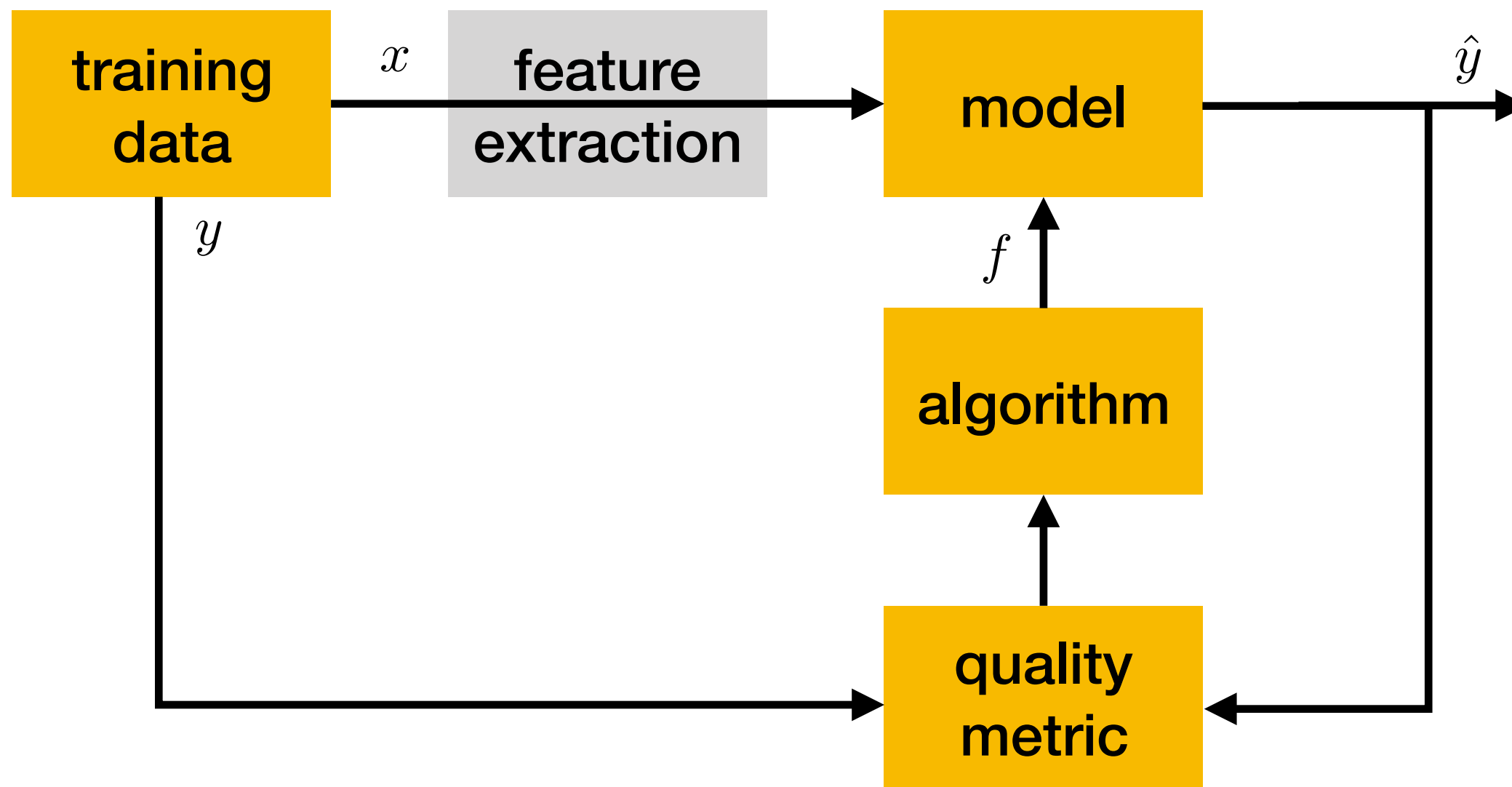
Commonly used loss functions for Classification

Recall **linear regression**

linear model: $f(x) = w_0 + w_1x[1] + w_2x[2] + \dots$

We choose square quality metric: $\text{MSE}(w) = \frac{1}{N} \sum_{i=1}^N (f(x_i) - y_i)^2$

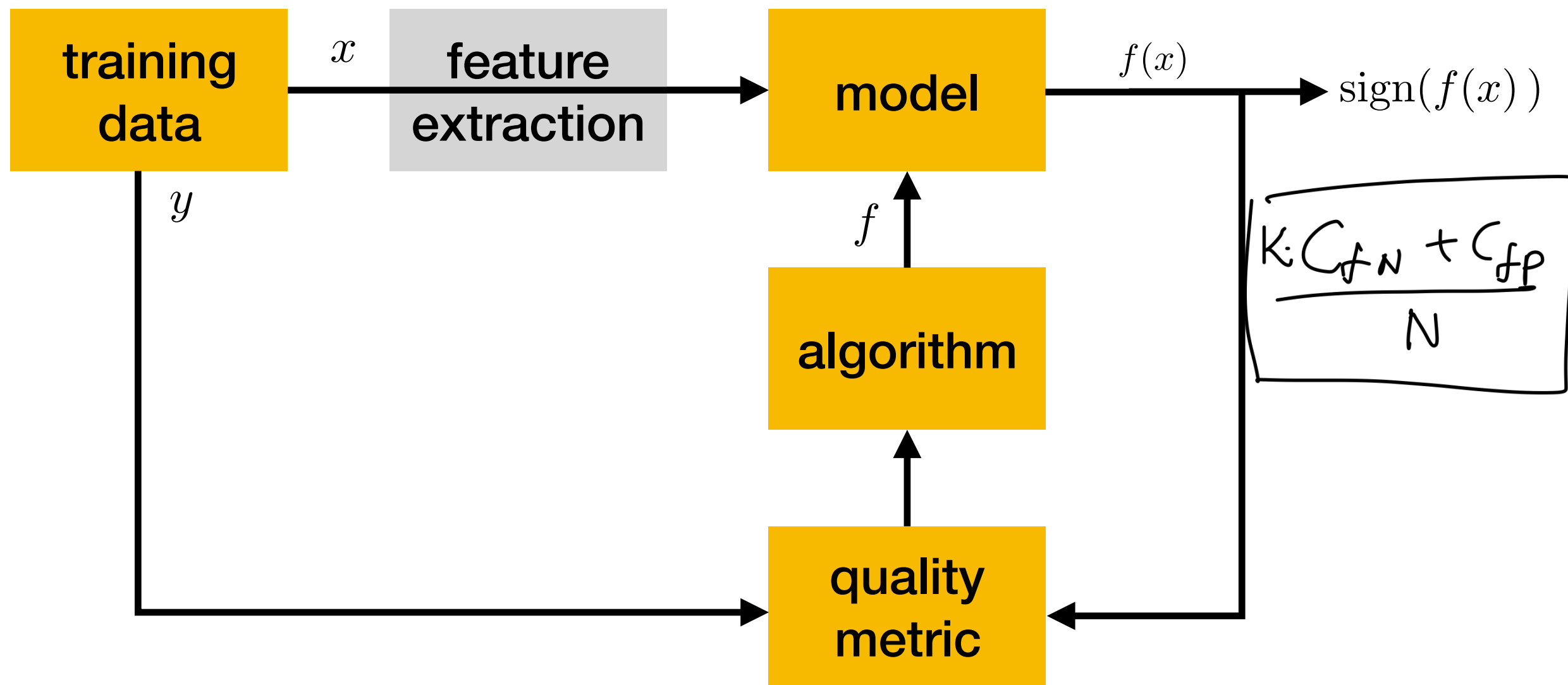
In training: find w that minimizes $\frac{1}{N} \sum_{i=1}^N (f(x_i) - y_i)^2$



Can we do the same for classification?

linear model: $f(x) = w_0 + w_1x[1] + w_2x[2] + \dots$

We choose k-Neyman-Pearson: $\text{Loss}(w) = \frac{1}{N} \sum_{i=1}^N \left\{ \underbrace{\kappa \mathbb{I}(\text{sign}(f(x_i)) = -1 \text{ and } y_i = 1)}_{\text{false negative}} + \underbrace{\mathbb{I}(\text{sign}(f(x_i)) = 1 \text{ and } y_i = -1)}_{\text{false positive}} \right\}$



Linear (Boolean) classifier

- You train a **linear model** of the form

$$f(x) = w_0 + w_1 h_1(x) + w_2 h_2(x) + \dots$$

- Prediction is

$$\hat{y} = \text{sign}(f(x))$$

- Ideally, we would like to find the weights \mathbf{w} , that minimizes **error rate** or more generally **Neyman-Pearson error**

$$\frac{\kappa C_{fn} + C_{fp}}{N} =$$

$$\frac{1}{N} \sum_{i=1}^N \left\{ \underbrace{\kappa \mathbb{I}(\text{sign}(f(x_i)) = -1 \text{ and } y_i = 1)}_{\text{false negative}} + \underbrace{\mathbb{I}(\text{sign}(f(x_i)) = 1 \text{ and } y_i = -1)}_{\text{false positive}} \right\}$$

Notations

- So far we used the notation

$$f(x) = w_0 + w_1 h_2(x) + w_2 h_x(x) + \dots$$

for the **continuous** linear model, and

$$\hat{y} = \text{sign}(f(x))$$

for the **discrete** prediction

- From now on, we will also use \hat{y} to denote the continuous valued model:

$$\hat{y} = f(x) = w_0 + w_1 h_2(x) + w_2 h_x(x) + \dots$$

to not introduce additional notation

- It should be clear from context which one we mean by \hat{y}

Training a linear classifier

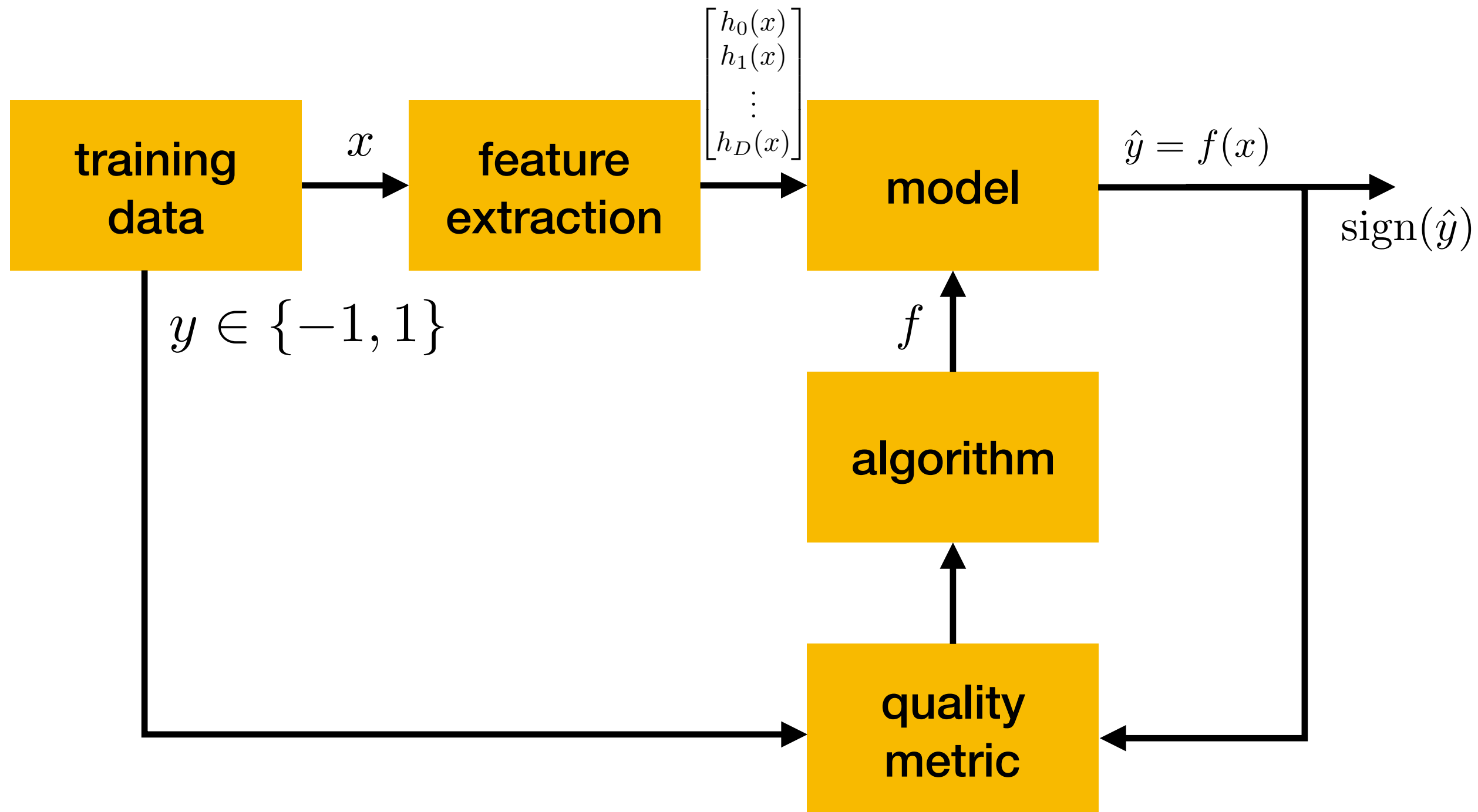
- Given a linear model

$$\hat{y} = f(x) = w_0 + w_1 h_1(x) + w_2 h_2(x) + \dots$$

- Neyman-Pearson error cannot be directly optimized (more on this later)
- Instead, in training classifiers, we minimize a loss of the form

$$\frac{1}{N} \sum_{i=1}^N \ell(\underbrace{\hat{y}_i}_{f(x_i) = w^T x_i = w_0 + w_1 x_1 [1] + \dots}, y_i)$$

- And find parameters \mathbf{w} , that minimize this loss function
- The choice of $\ell(\hat{y}, y)$ depends on the application
- One can add regularization: $\lambda r(w)$



- Recipe for training classifiers
 - 1. Train a continuous valued model, as if regression but with special choices of the loss
 - 2. For prediction take $\text{sign}(f(x))$
 - 3. The score $f(x)$ tells us how confident we are in the prediction

Loss function for Boolean classification

- We need to design loss function $\ell(\hat{y}, y)$
- Note that $w^T x = w_0 + w_1 x[1] + w_2 x[2] + \dots$
 - prediction $\hat{y} = w^T x$ can take any values
 - but y can only take $+1$ or -1
- So in order to specify $\ell(\hat{y}, y)$, we only need to give two functions (of scalar \hat{y})
 - $\ell(\hat{y}, -1)$ is how much \hat{y} irritates us when $y = -1$
 - $\ell(\hat{y}, 1)$ is how much \hat{y} irritates us when $y = 1$
- typically, one chooses those two functions to be symmetric, but appropriately scaled to reflect that false negatives irritates us factor κ more than false positives:

$$\ell(\hat{y}, 1) = \kappa \ell(-\hat{y}, -1)$$

Neyman-Pearson loss

$$\sum_i \ell(f(x_i), y_i)$$

$$\tilde{\ell}(y_i, f(x_i))$$

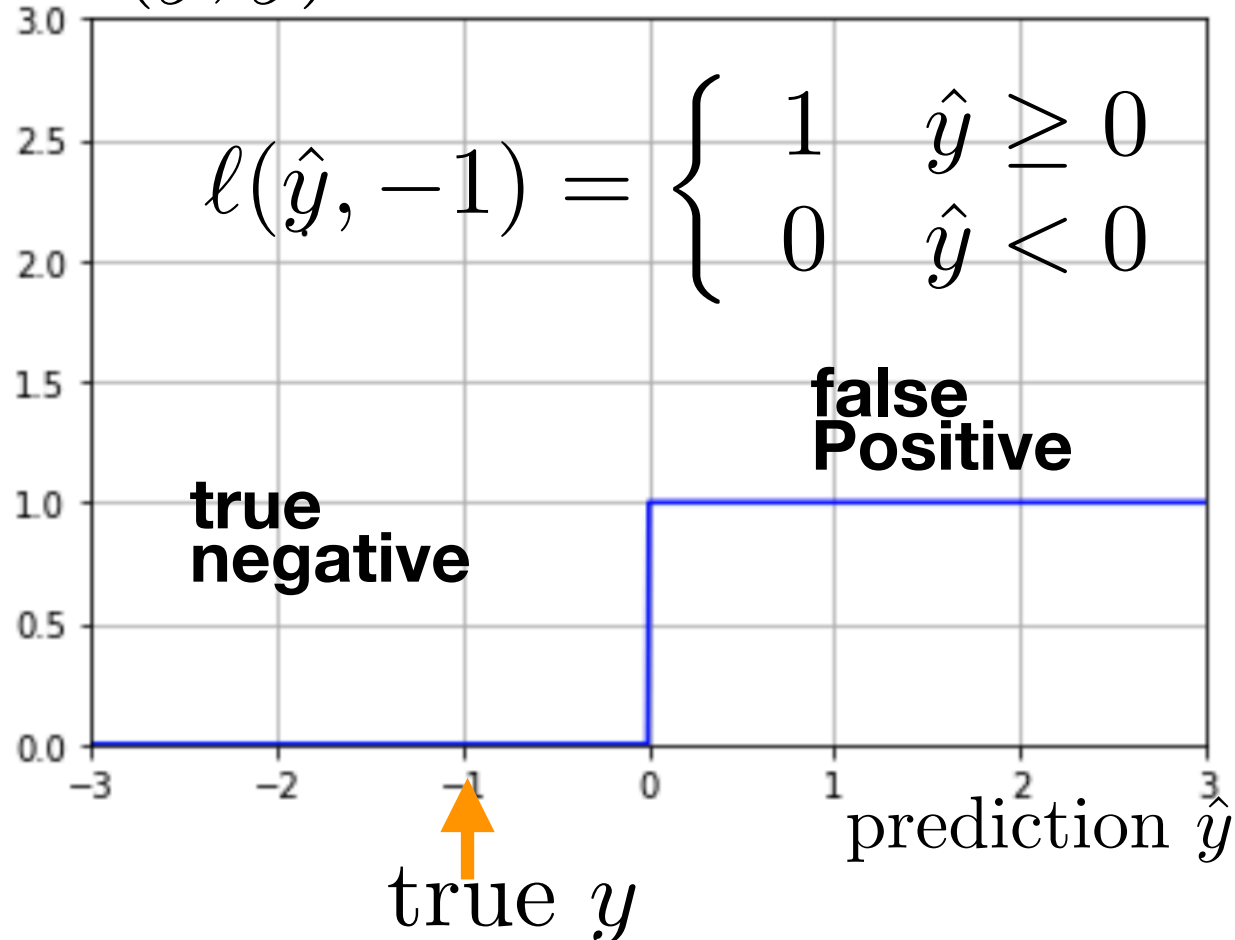
- Neyman-Pearson loss is

$$\ell(\hat{y}, -1) = \begin{cases} 1 & \hat{y} \geq 0 \\ 0 & \hat{y} < 0 \end{cases}$$

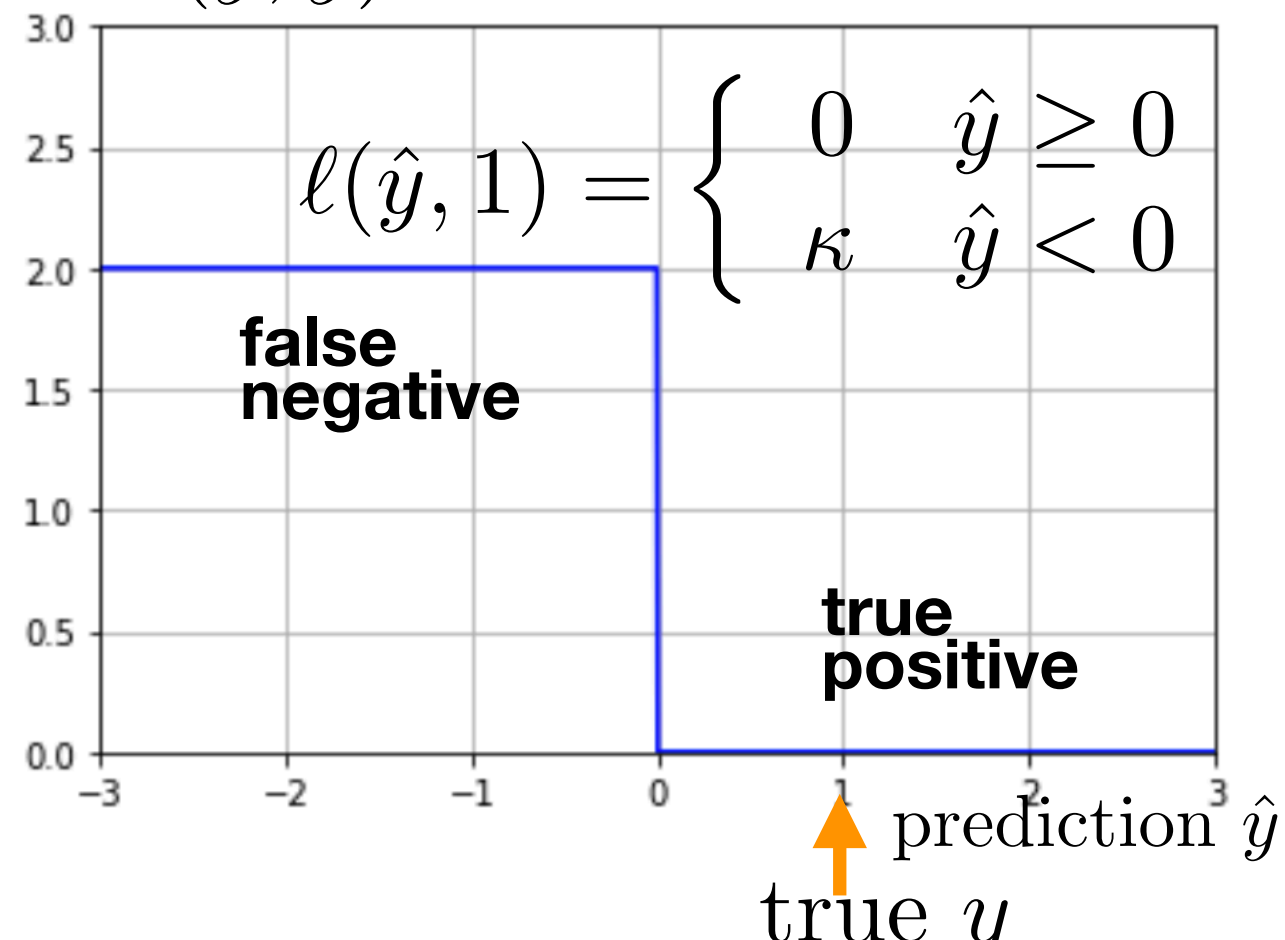
$$\ell(\hat{y}, 1) = \begin{cases} 0 & \hat{y} \geq 0 \\ \kappa & \hat{y} < 0 \end{cases}$$

- Neyman-Pearson loss computed on the training data is (training) Neyman-Pearson error

loss $\ell(\hat{y}, y)$



loss $\ell(\hat{y}, y)$



Problem with Neyman-Pearson loss

- Neyman-Pearson loss is not differentiable, or even continuous (And certainly not convex)
- Its gradient is zero or does not exist
- Gradient based optimizer does not know how to improve the model

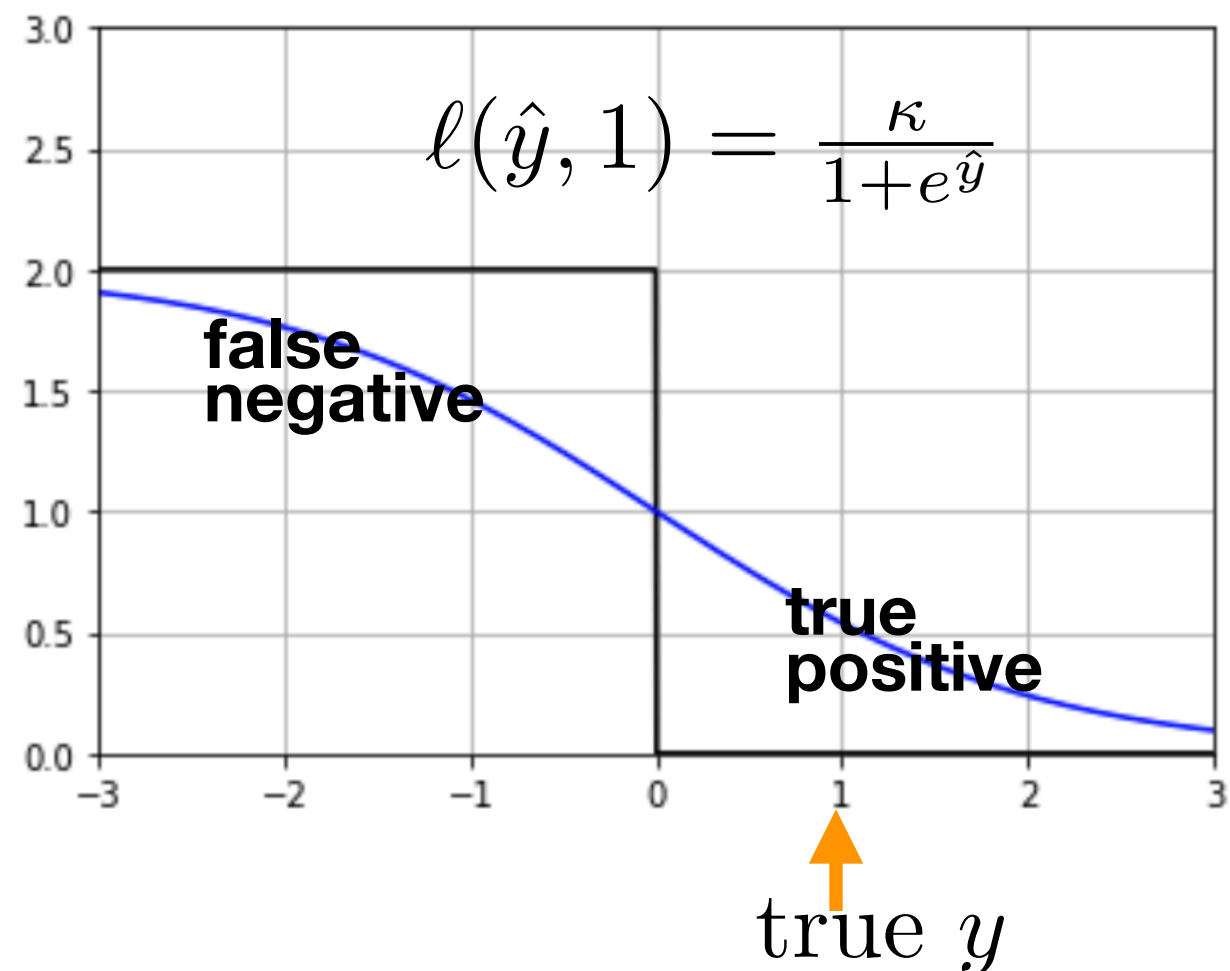
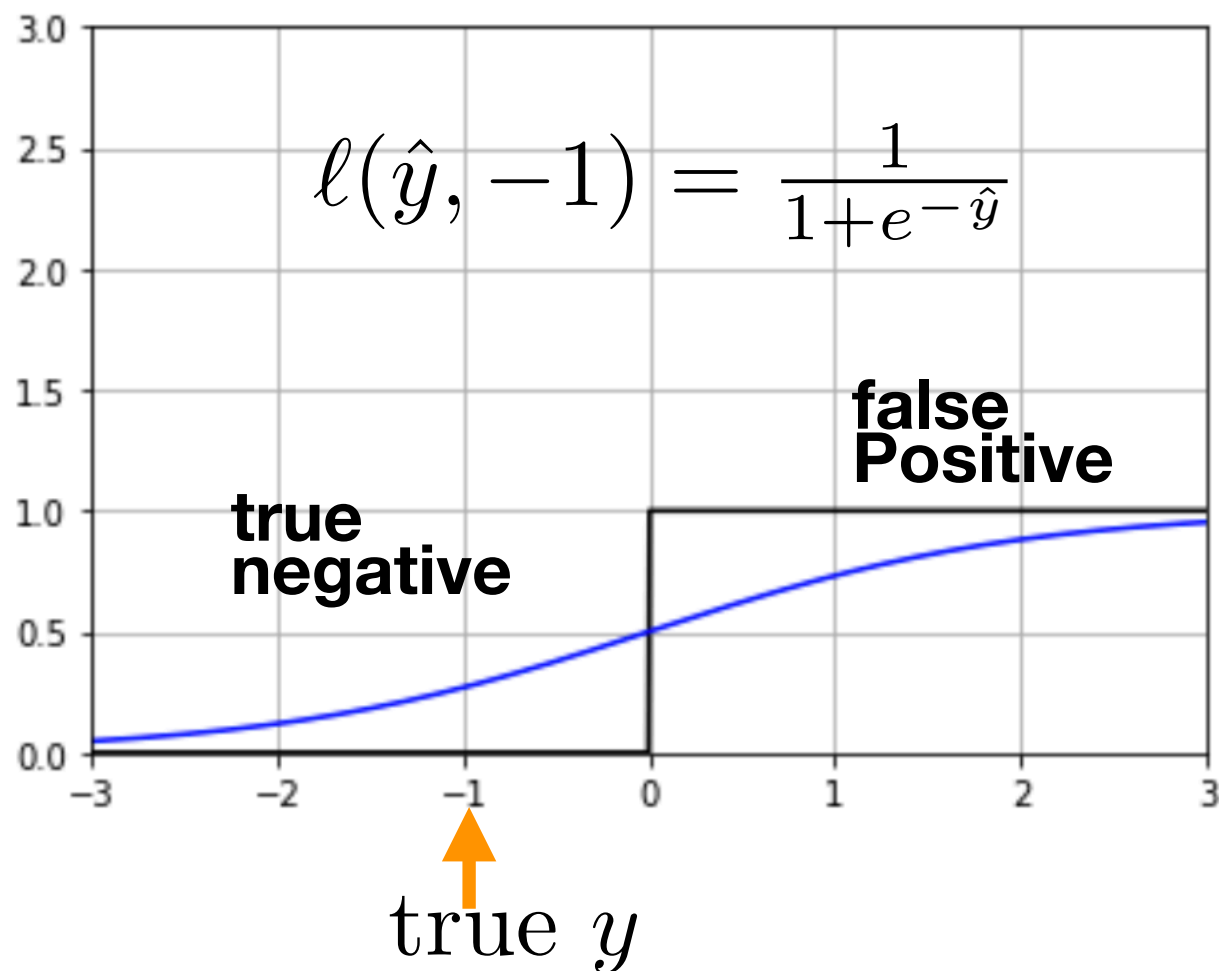
Ideas of proxy loss

- We get better results using proxy losses that
 - approximate, or captures the flavor of, the Neyman-Pearson loss
 - Is more easily optimized (e.g. convex or non-zero derivatives)
- concretely, we want **proxy loss function**
 - with $\ell(\hat{y}, -1)$ small when $\hat{y} < 0$ and larger when $\hat{y} > 0$
 - with $\ell(\hat{y}, 1)$ small when $\hat{y} > 0$ and larger when $\hat{y} < 0$
 - Which has other nice characteristics, e.g., differentiable or convex

Sigmoid loss

- Differentiable approximation of Neyman-Pearson loss
- But not convex in \hat{y}
- The two losses sum to one, if $k=1$
- Softer (or smoothed) version of the N-P loss

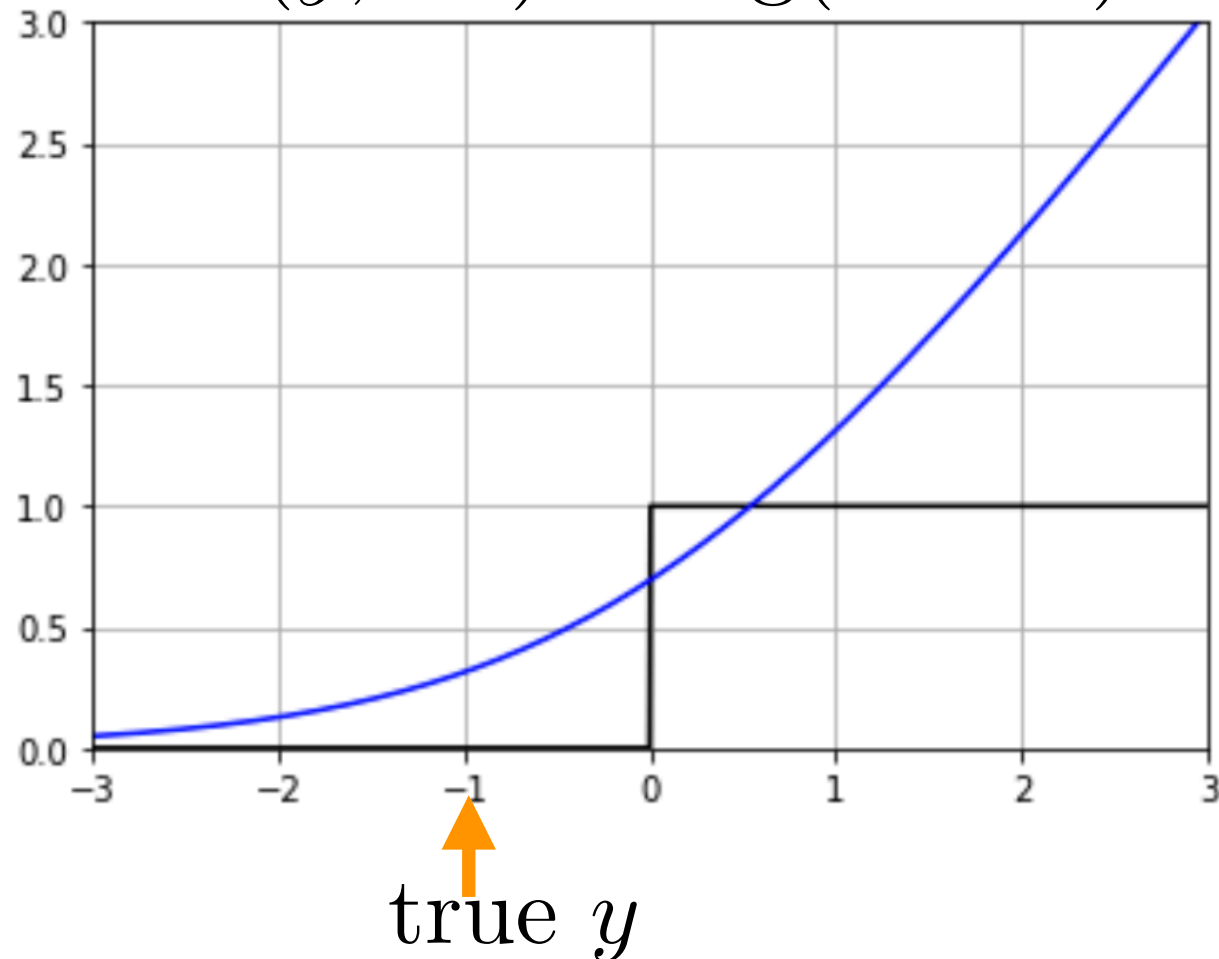
$$\frac{1}{1+e^{-y}} + \frac{1}{1+e^y} = \frac{e^y}{e^y+1} + \frac{1}{1+e^y} = 1$$



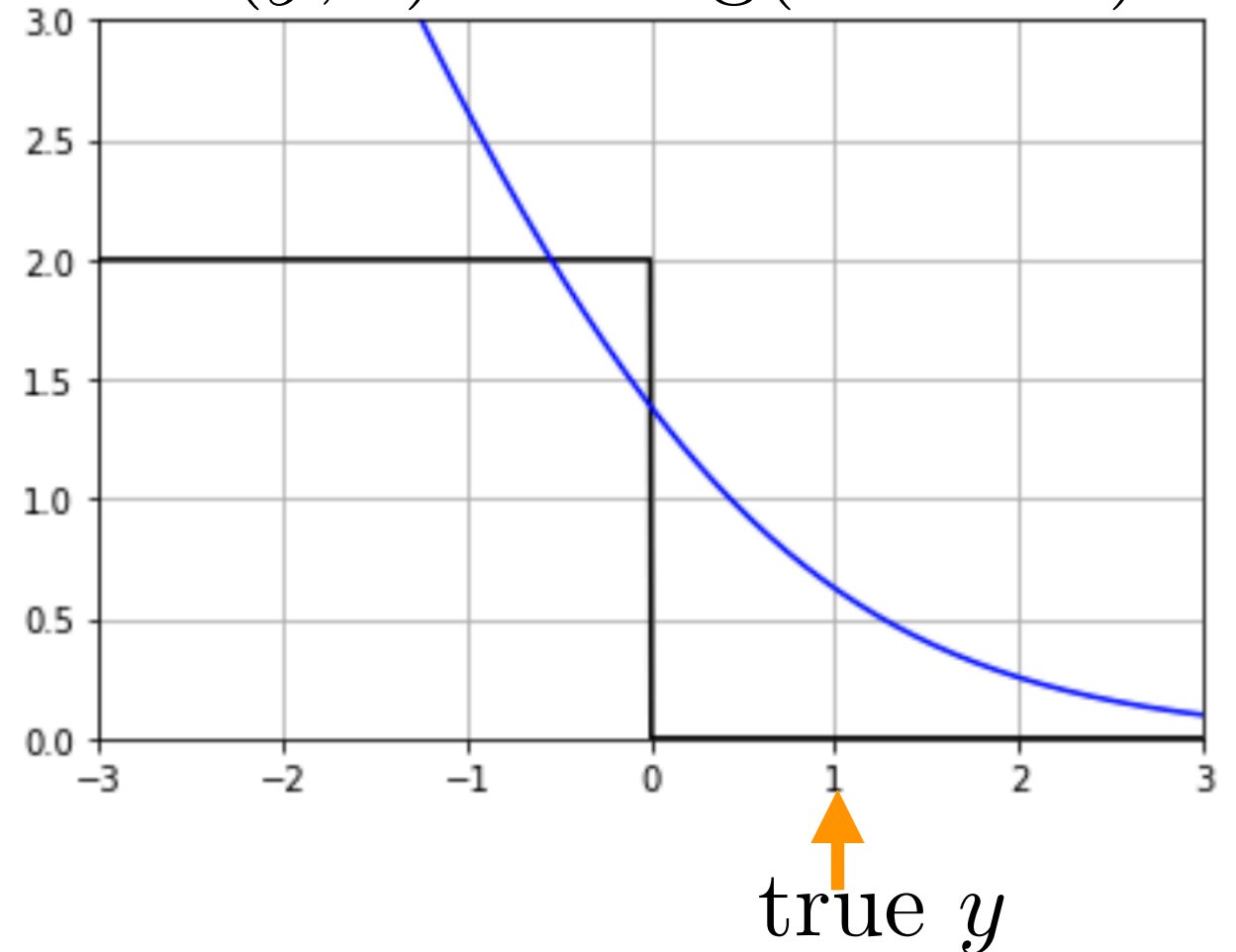
Logistic loss

- Differentiable and convex in \hat{y}
- approximation of Neyman-Pearson

$$\ell(\hat{y}, -1) = \log(1 + e^{\hat{y}})$$



$$\ell(\hat{y}, 1) = \kappa \log(1 + e^{-\hat{y}})$$

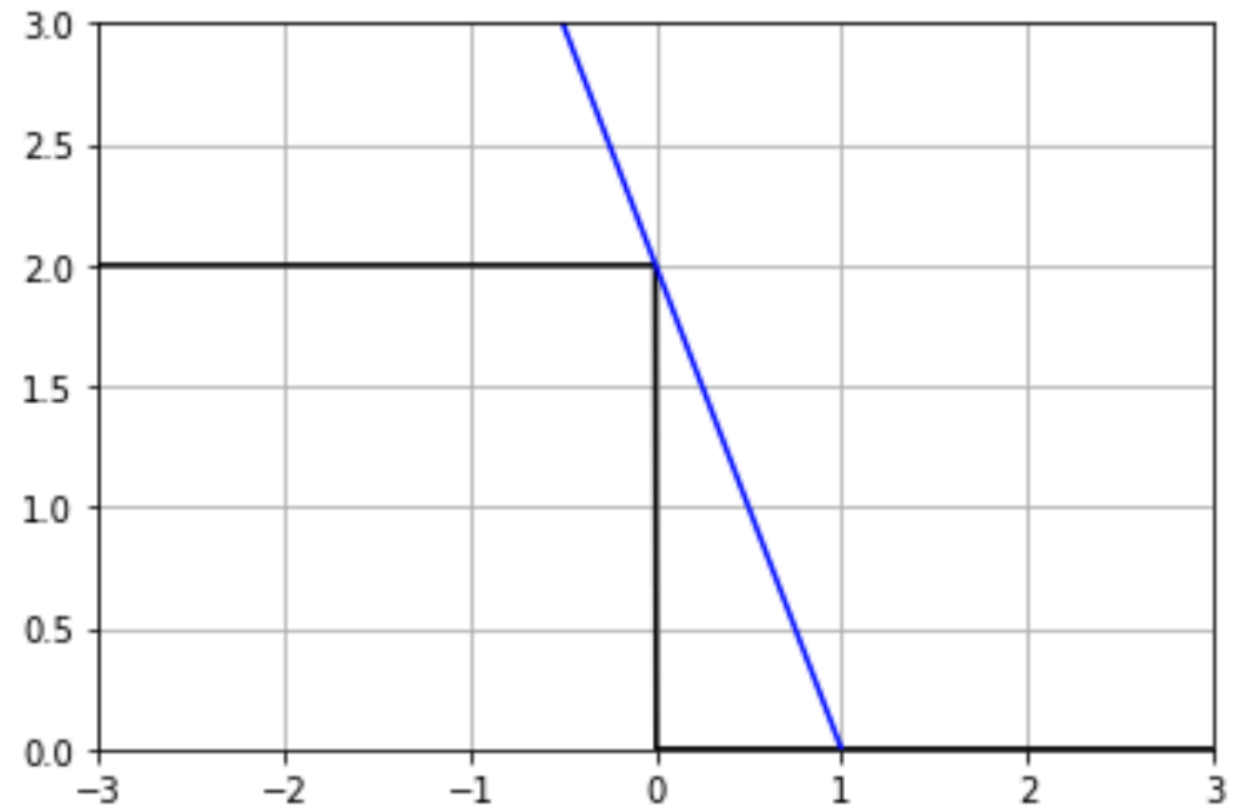
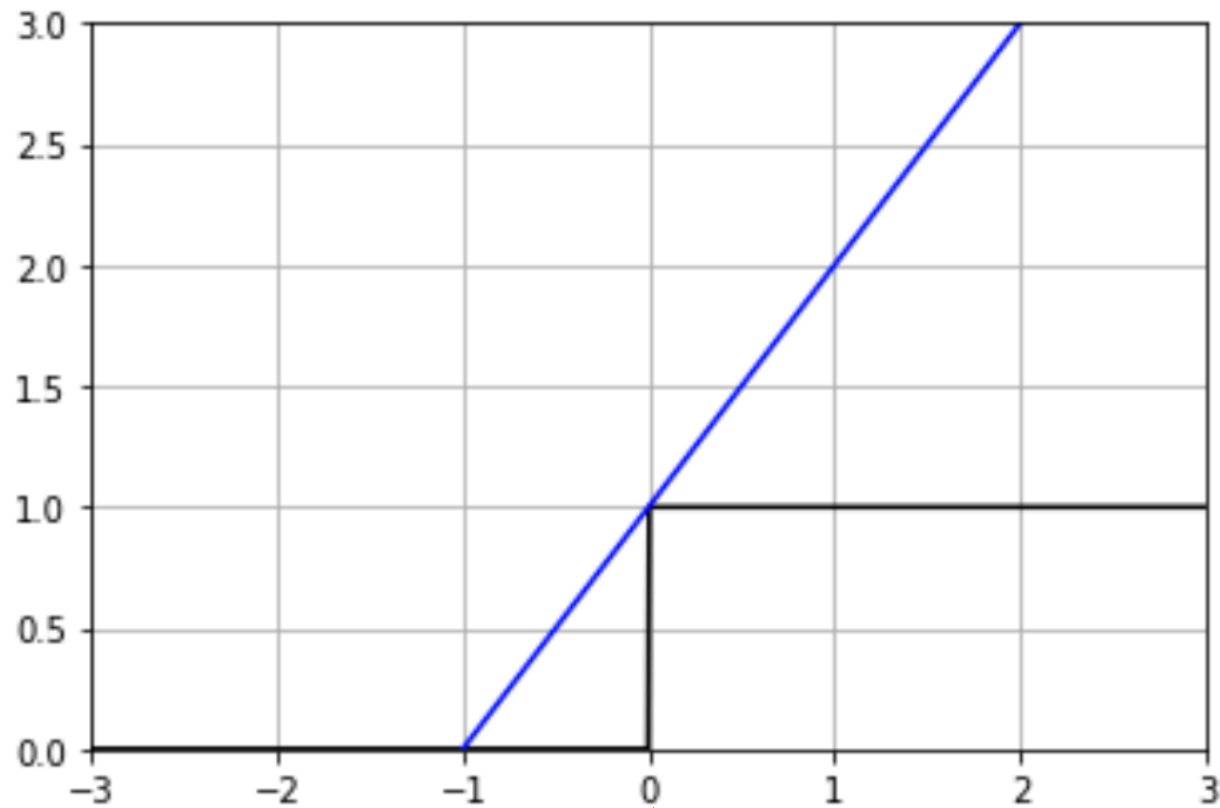


Hinge loss

- Non-differentiable but convex approximation of Neyman-Pearson loss

$$\ell(\hat{y}, -1) = [1 + \hat{y}]^+$$

$$\ell(\hat{y}, 1) = \kappa[1 - \hat{y}]^+$$



true y

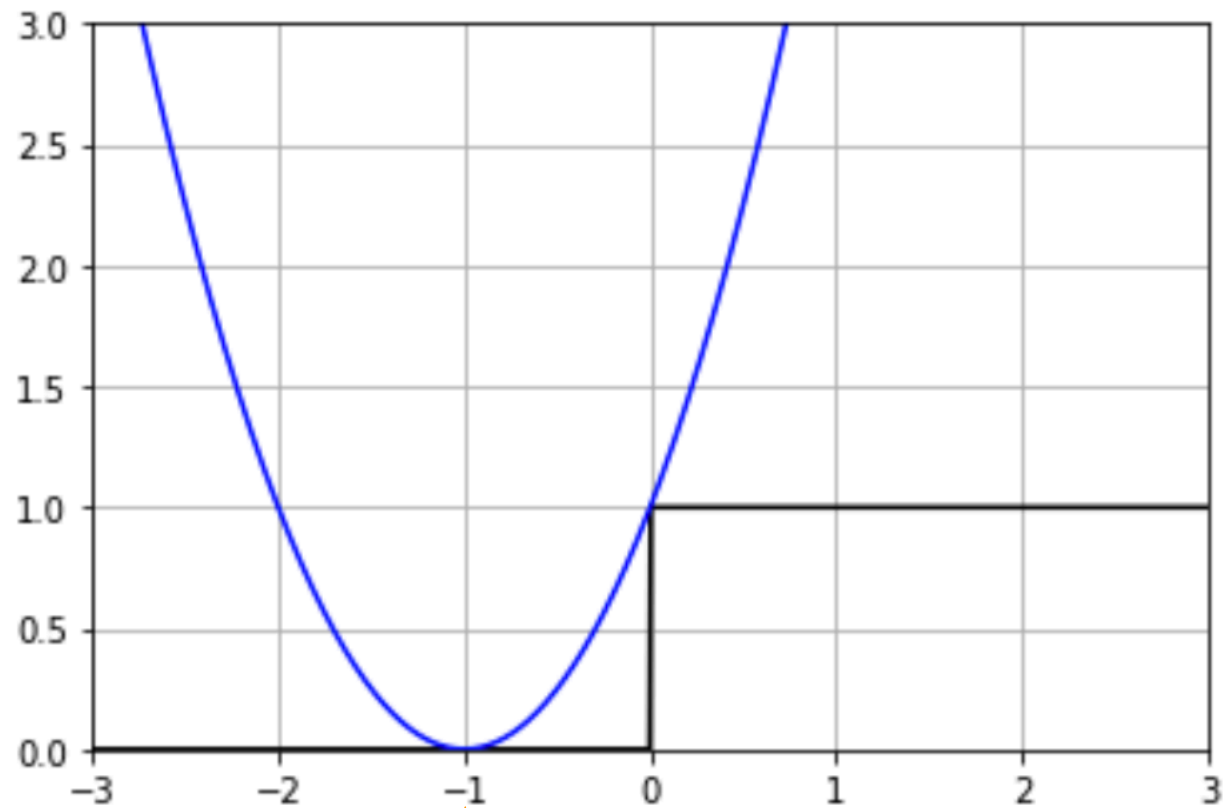
where $[x]^+ = \max\{0, x\}$

true y

Square loss

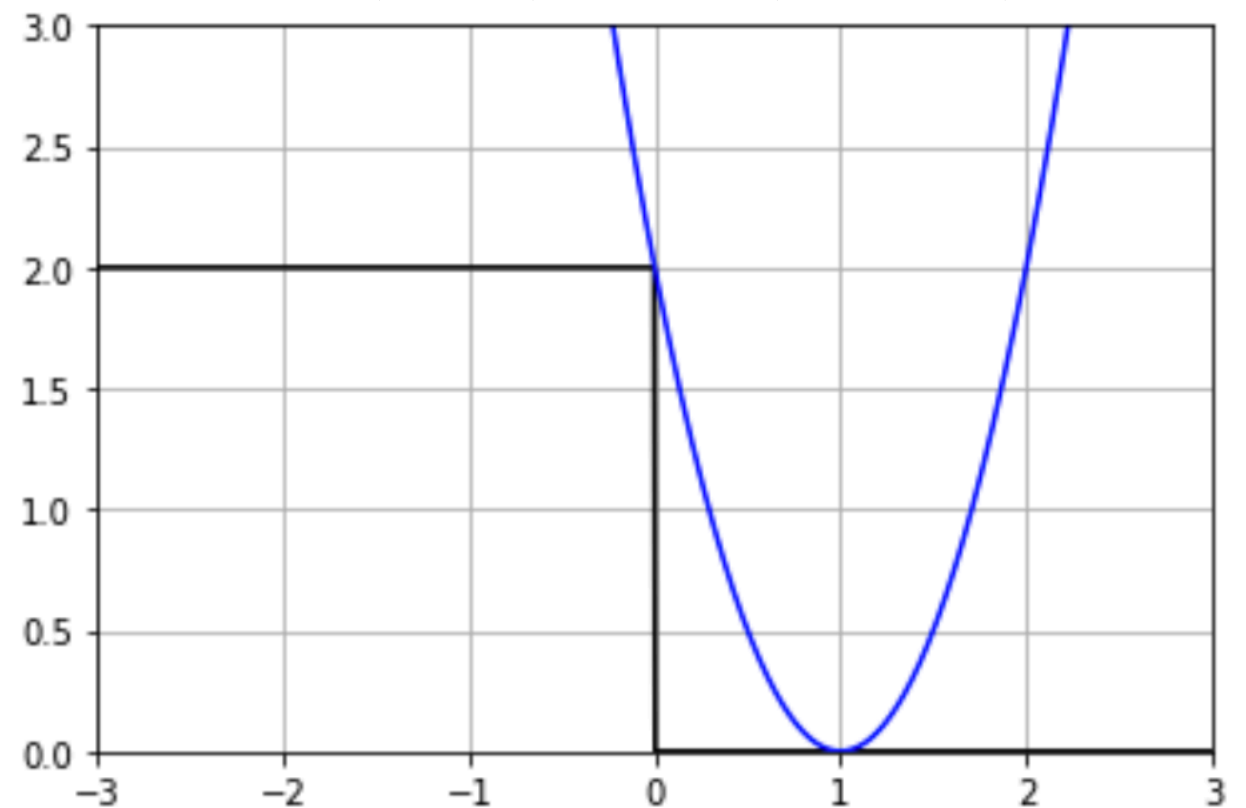
- Not only is it convex, square loss is easy to minimize (has a closed form solution)

$$\ell(\hat{y}, -1) = (1 + \hat{y})^2$$

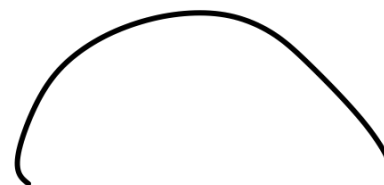
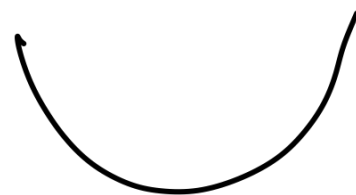


↑
true y

$$\ell(\hat{y}, 1) = \kappa(1 - \hat{y})^2$$



↑
true y



Commonly used Boolean classifiers

Squared loss classifier

- Uses **sum of squares loss** (a.k.a. L2 loss, Mean Squared Error (MSE), Residual Sum of Squares (RSS))

$$\text{minimize}_w \mathcal{L}(w) = \frac{1}{N} \left(\sum_{i:y_i=-1} (1 + \hat{y}_i)^2 + \kappa \sum_{i:y_i=1} (1 - \hat{y}_i)^2 \right)$$

together with a choice of your regularizer

- This is particularly easy to optimize, if the regularizer is also L2 regularizer

Logistic regression

- Uses **logistic loss**

$$\text{minimize}_w \mathcal{L}(w) = \frac{1}{N} \left(\sum_{i:y_i=-1} \log(1 + e^{\hat{y}_i}) + \kappa \sum_{i:y_i=1} \log(1 + e^{-\hat{y}_i}) \right)$$

with a choice of a regularizer

- Is a convex optimization if the regularizer is convex, and the minimizer can be found efficiently

)

Support vector machine (SVM)

- Uses hinge loss

$$\text{minimize}_w \mathcal{L}(w) = \frac{1}{N} \left(\sum_{i:y_i=-1} [1 + \hat{y}_i]^+ + \kappa \sum_{i:y_i=1} [1 - \hat{y}_i]^+ \right)$$

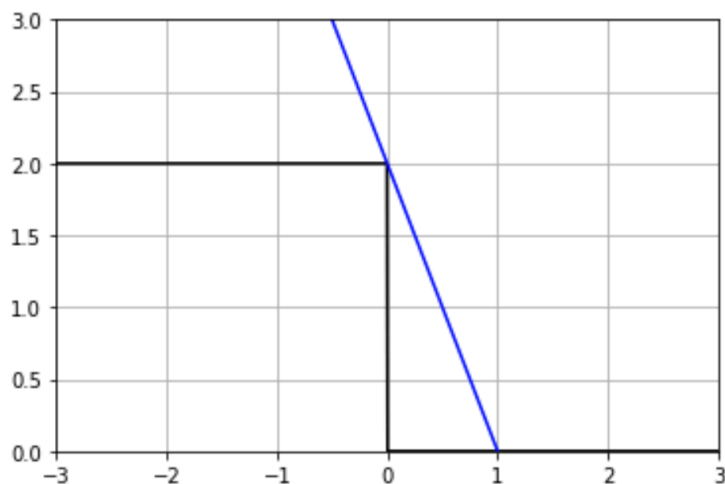
with sum of squares regularizer

where $[x]^+ = \max\{0, x\}$

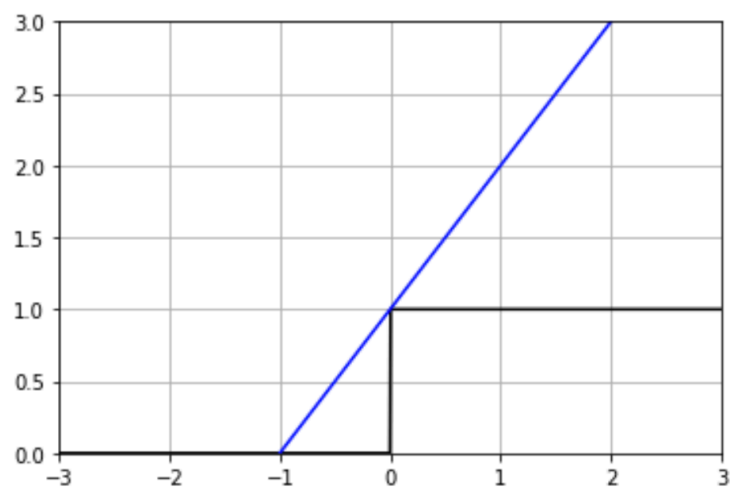
- It is a convex minimization

Support vector machine (SVM)

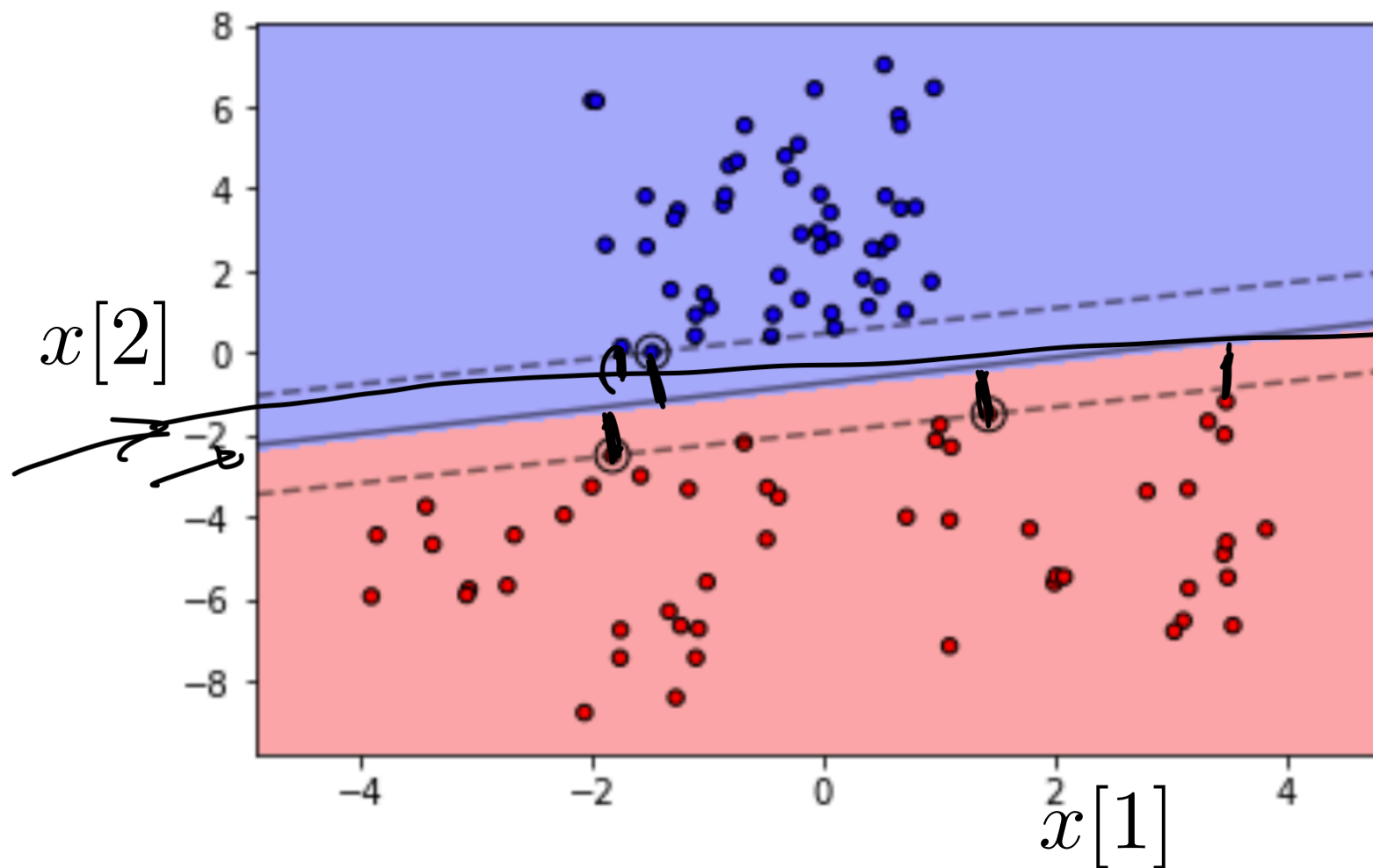
$$\ell(\hat{y}, 1) = \kappa[1 - \hat{y}]^+$$



$$\ell(\hat{y}, -1) = [1 + \hat{y}]^+$$



- Linear model is trained on the hinge loss shown on the left with $k=1$
- Resulting prediction is shown below



- As we predict with $\text{sign}(\hat{y})$, the decision boundary is at
$$\hat{y} = w_0 + w_1x[1] + w_2x[2] = 0$$
- Dotted lines show the points where $\hat{y} = w_0 + w_1x[1] + w_2x[2] = \pm 1$
- What is the training error?

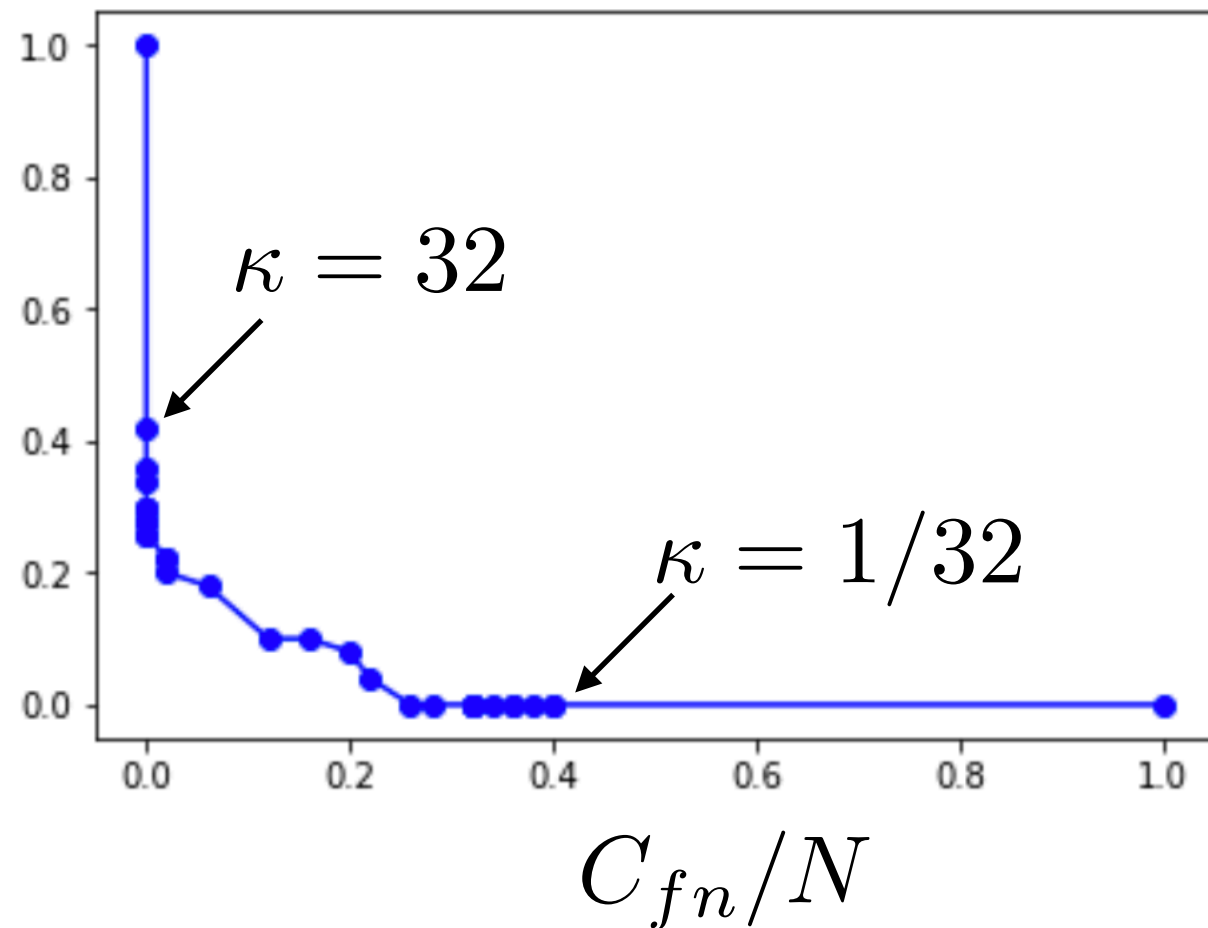
Receiver Operating Characteristic (ROC)

Receiver Operating Characteristic (ROC)

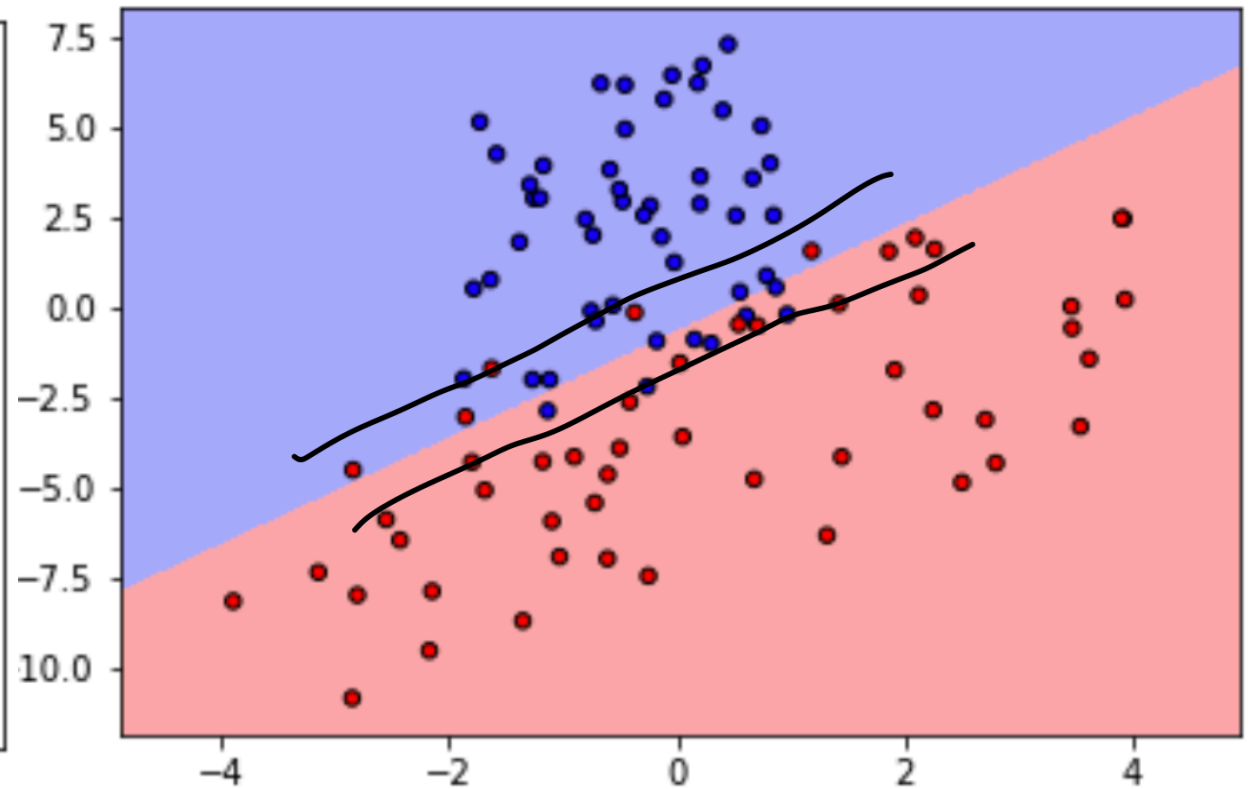
- Always abbreviated as **ROC**, comes from WWII
- Explores the tradeoff between false negative and false positive rates
- Typical recipe for evaluating performance of a classifier
 - 1. Construct a classifier for many values of k
 - For each k , select the regularization hyper-parameter via cross-validation, that minimizes Neyman-Pearson loss on test data set
 - 2. Plot the computed pair (false negative rate, false positive rate) on a 2-D plot.
- Connecting all the dots gives you **ROC curve** (when viewed upside-down)

Example: ROC curve

C_{fp}/N



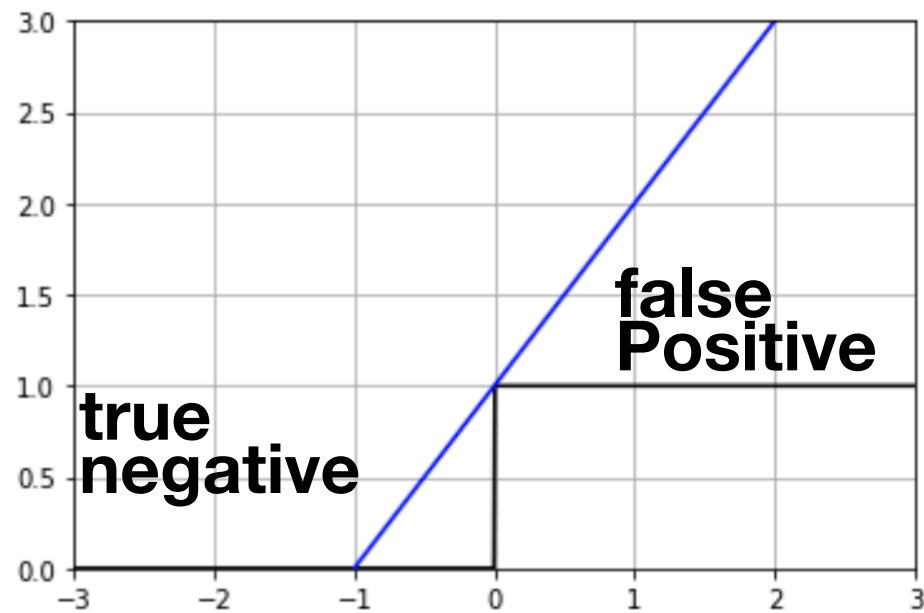
$\kappa = 0$



- SVM with various k
- Left hand plot shows training error pairs $(C_{fn}/N, C_{fp}/N)$
- Right hand plot shows minimum error classifier (i.e. $k=1$)

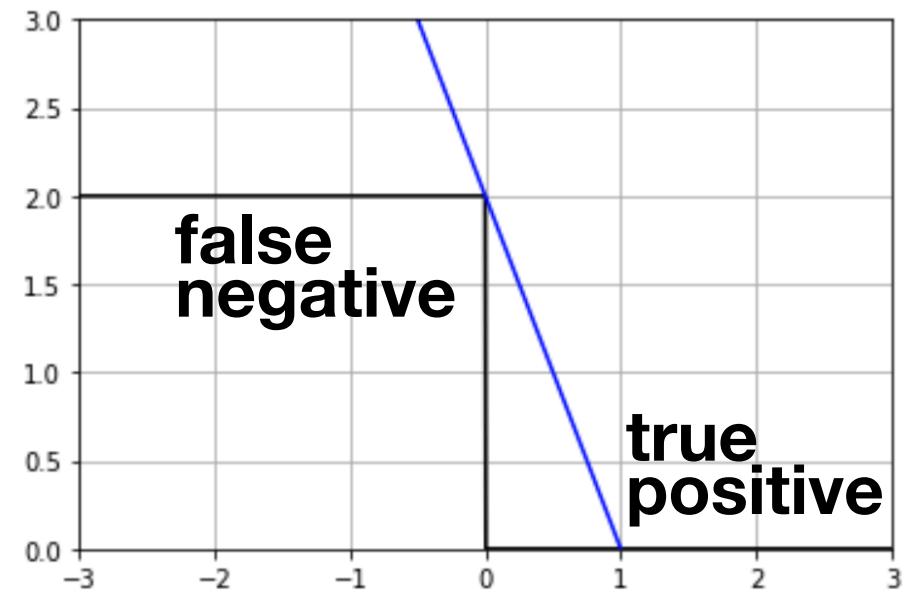
Example

$$\ell(\hat{y}, -1) = [1 + \hat{y}]^+$$

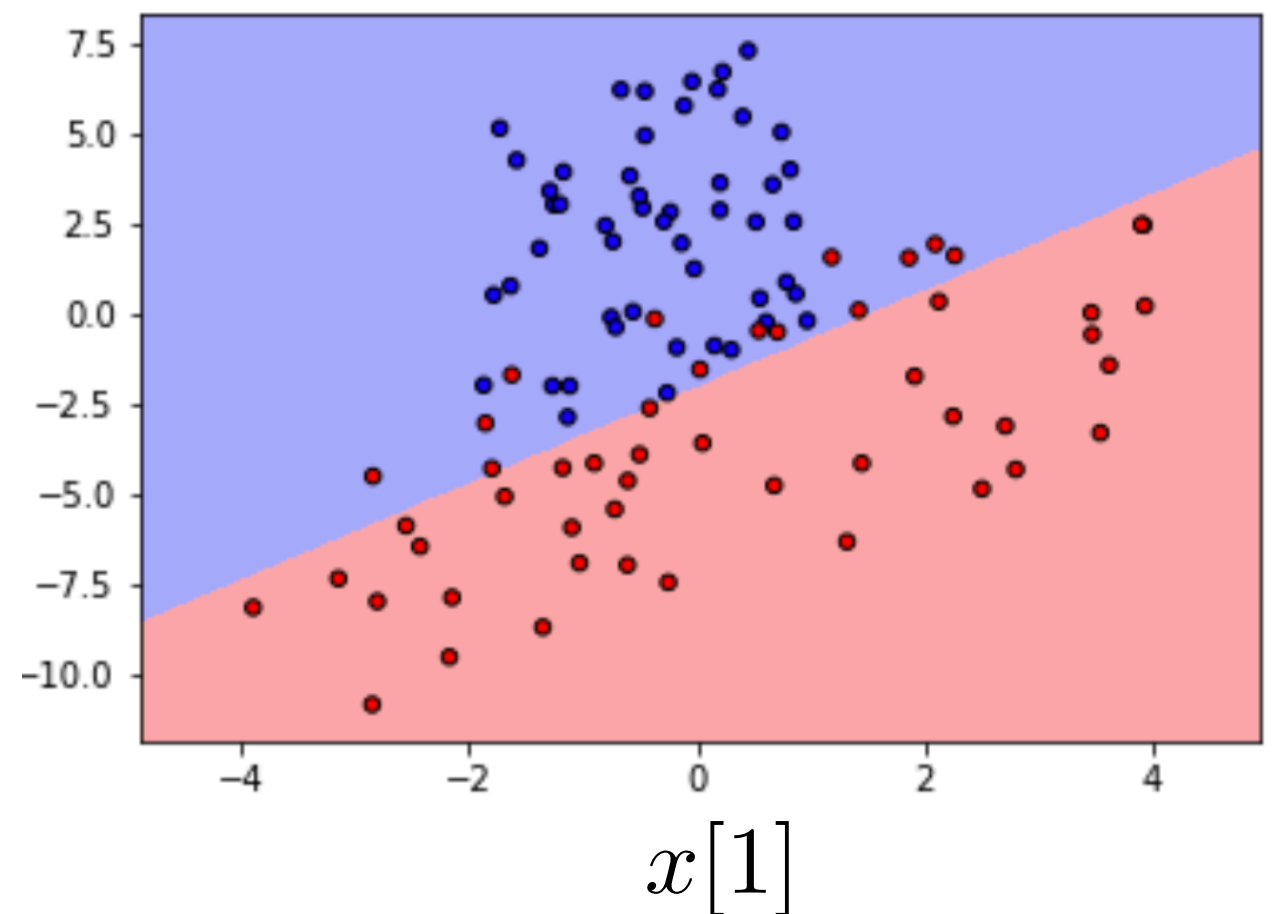
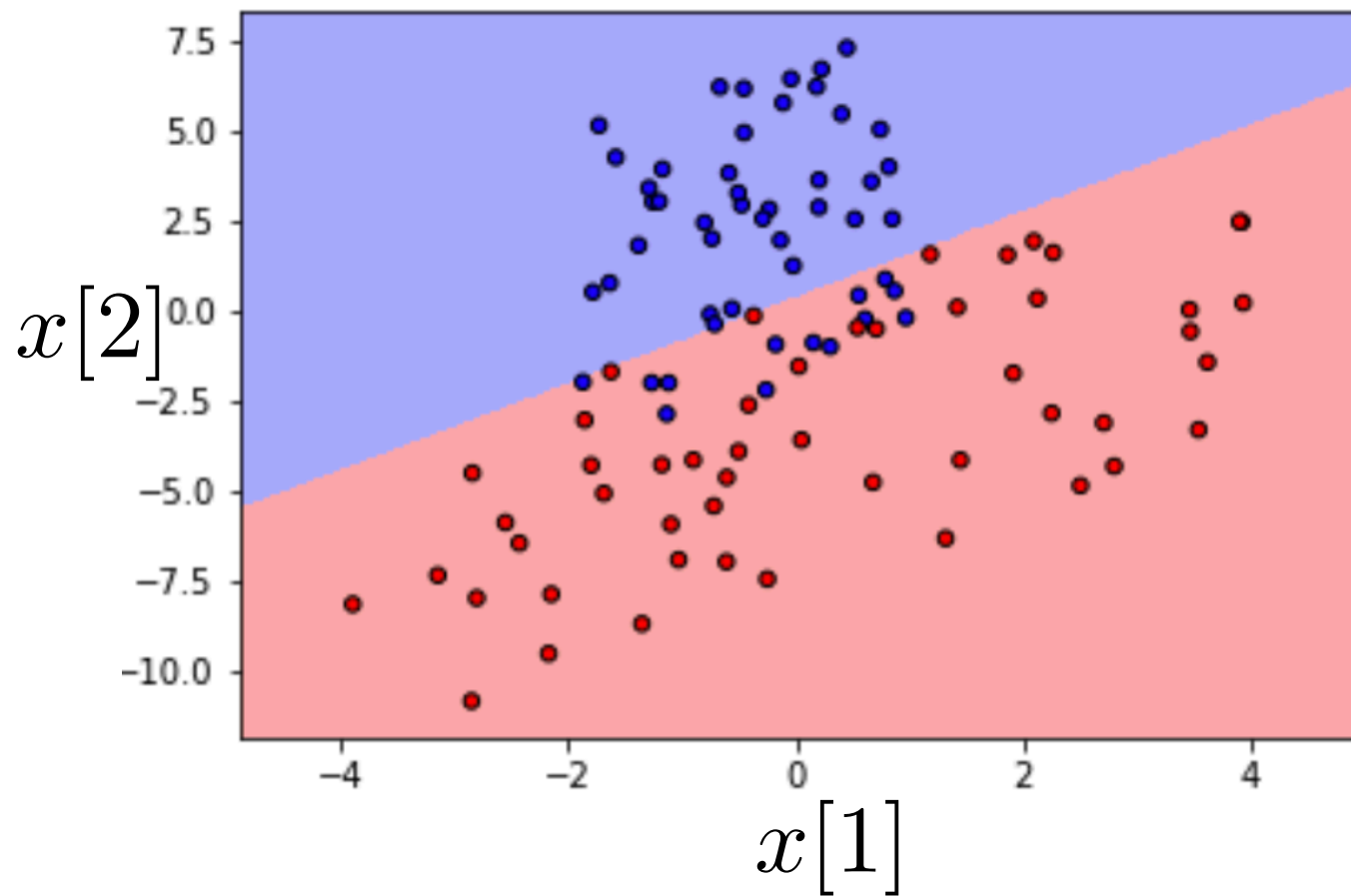


$$\kappa = 0.4$$

$$\ell(\hat{y}, 1) = \kappa[1 - \hat{y}]^+$$

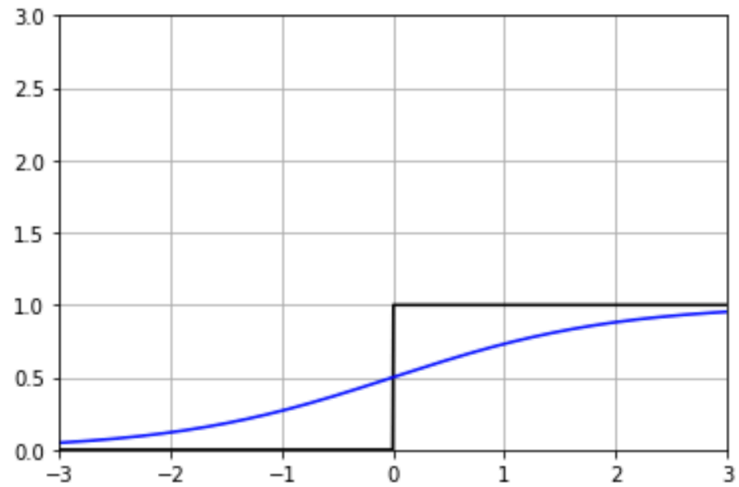


$$\kappa = 4$$



Probabilistic interpretation of **logistic regression**

- When $\kappa = 1$, we get the following losses for each data point x



$$\left(\underbrace{\frac{1}{1 + e^{-w^T x}}}_{\ell(\hat{y}, -1)}, \underbrace{\frac{1}{1 + e^{w^T x}}}_{\ell(\hat{y}, 1)} \right)$$

when using **sigmoid loss** that is trained with a linear model

- They are
 - Non-negative
 - sum to one, and
 - they measure how likely it is that the point x has label +1 (or -1) respectively
- One can view it as an estimation of the probability

$$\left(\mathbb{P}(y_i = +1 | x), \mathbb{P}(y_i = -1 | x) \right)$$

•

Probabilistic interpretation of **logistic regression**

- Then taking the sign of the linear predictor to make final decision is simply taking a label that is more likely:

$$\hat{y} = \text{sign}(w^T x) \quad \iff \quad \hat{y} = \begin{cases} +1 & \frac{1}{1+e^{-w^T x}} > \frac{1}{1+e^{w^T x}} \\ -1 & \text{otherwise} \end{cases}$$

- and **logistic regression** can be interpreted as **Maximum Likelihood Estimator** under the probabilistic model with sigmoid function:

- $$\left(\underbrace{\frac{1}{1+e^{-w^T x}}}_{\mathbb{P}(y_i=+1|x_i)}, \underbrace{\frac{1}{1+e^{w^T x}}}_{\mathbb{P}(y_i=-1|x_i)} \right)$$

Maximum Likelihood Estimator (MLE)

- model: $\left(\underbrace{\frac{1}{1 + e^{-w^T x}}}_{\mathbb{P}(y_i = +1 | x_i)}, \underbrace{\frac{1}{1 + e^{w^T x}}}_{\mathbb{P}(y_i = -1 | x_i)} \right)$
- log-likelihood on a data point (x_i, y_i) : $\hat{y}_i = w_0 + w_1 x_i [1] + w_2 x_i [2]$

$$\text{log-likelihood} = \log \left(\mathbb{P}(y_i | x_i) \right) = \begin{cases} \log \left(\frac{1}{1 + e^{-w^T x_i}} \right) & \text{if } y_i = +1 \\ \log \left(\frac{1}{1 + e^{w^T x_i}} \right) & \text{if } y_i = -1 \end{cases}$$

- Maximum Likelihood Estimator is the one that maximizes the sum of all likelihoods on training data points

$$\text{maximize}_w \sum_{i: y_i = -1} \log \left(\frac{1}{1 + e^{\hat{y}_i}} \right) + \sum_{i: y_i = 1} \log \left(\frac{1}{1 + e^{-\hat{y}_i}} \right)$$

$\hat{y}_i = f(x_i)$

- Notice that this is exactly the **logistic regression** without any regularizers and with $k=1$

$$\text{minimize}_w \mathcal{L}(w) = \frac{1}{N} \left(\sum_{i: y_i = -1} \log(1 + e^{\hat{y}_i}) + \sum_{i: y_i = 1} \log(1 + e^{-\hat{y}_i}) \right)$$

