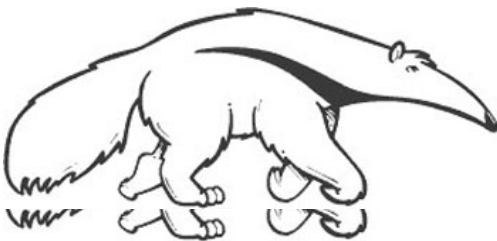


+

# Machine Learning and Data Mining

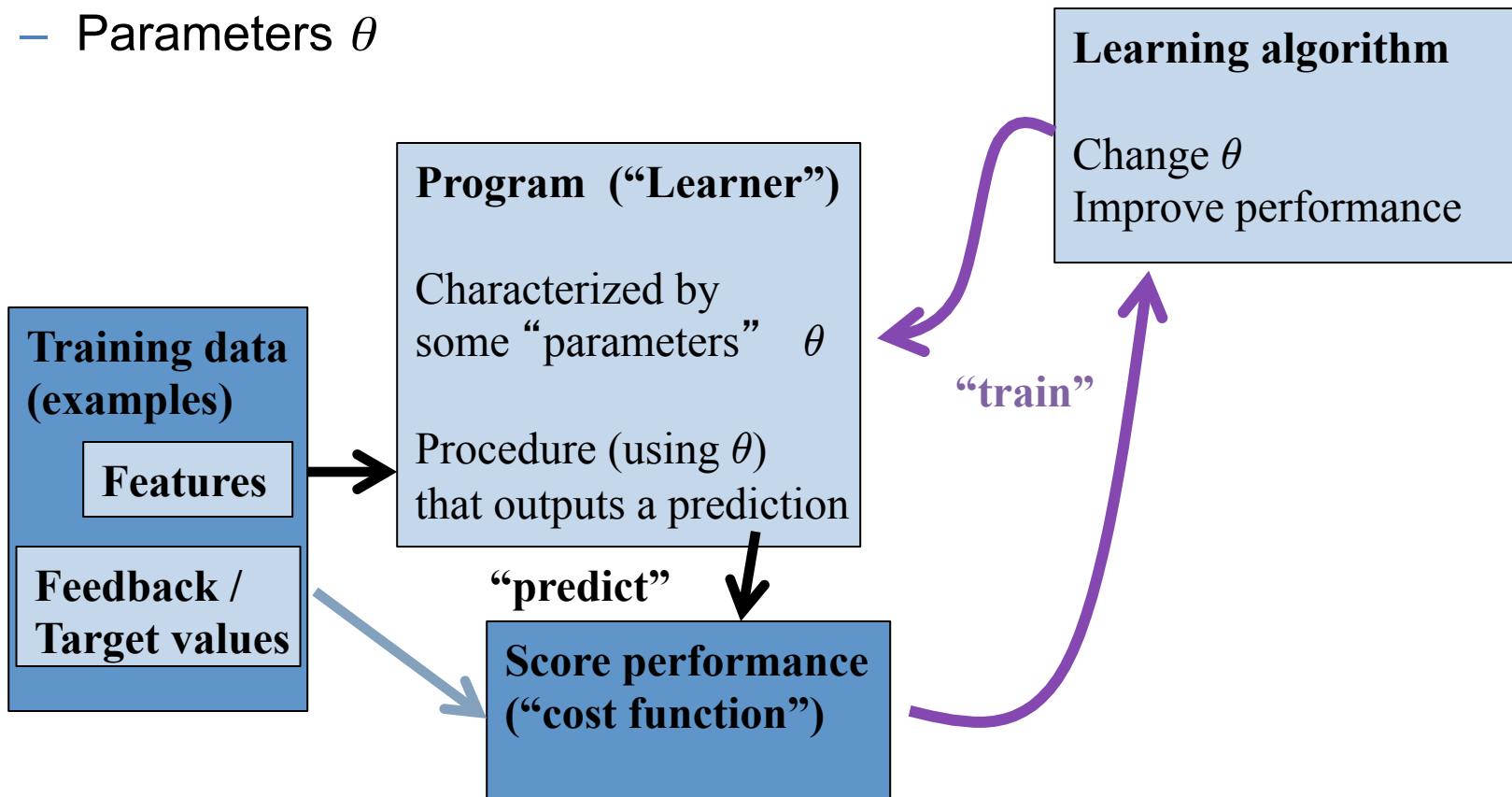
## Linear classification

Prof. Alexander Ihler

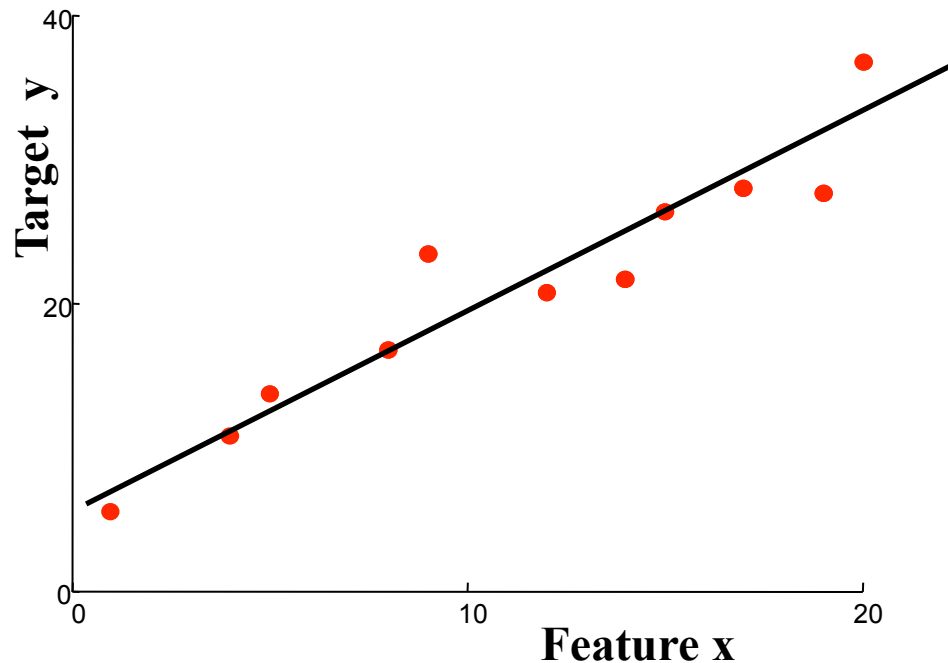


# Supervised learning

- Notation
  - Features  $x$
  - Targets  $y$
  - Predictions  $\hat{y} = f(x ; \theta)$
  - Parameters  $\theta$



# Linear regression



**“Predictor”:**

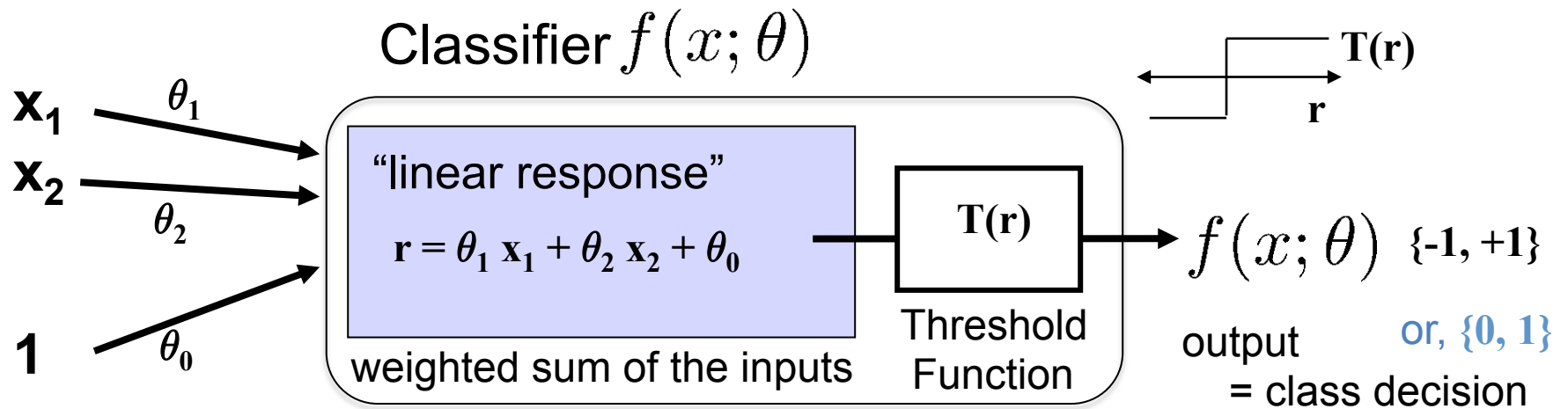
Evaluate line:

$$r = \theta_0 + \theta_1 x_1$$

return r

- Contrast with classification
  - Classify: predict discrete-valued target y
  - Initially: “classic” binary  $\{-1, +1\}$  classes; generalize later

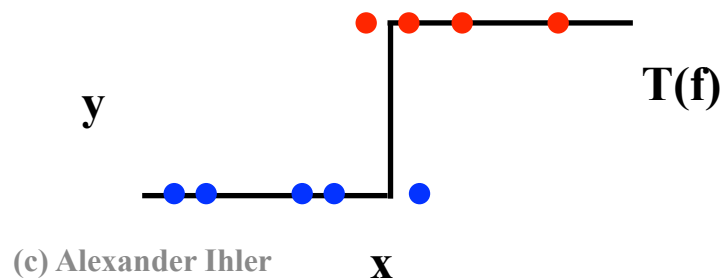
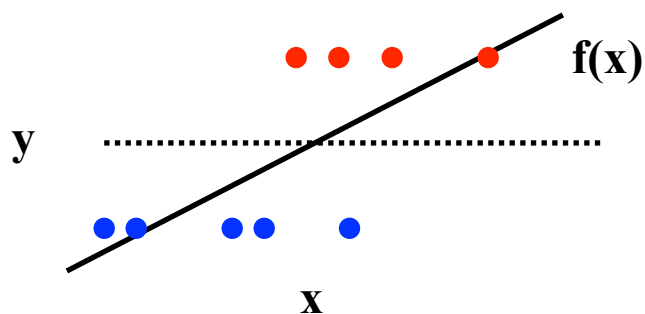
# Perceptron Classifier (2 features)



```

r = X.dot( theta.T );           # compute linear response
Yhat = (r > 0)                  # predict class 1 vs 0
Yhat = 2*(r > 0)-1              # or "sign": predict +1 / -1
# Note: typically convert classes to "canonical" values 0,1,...
# then convert back ("learner.classes[c]") after prediction
    
```

Visualizing for one feature “x”:



(c) Alexander Ihler

# Perceptrons

- Perceptron = a linear classifier
  - The parameters  $\theta$  are sometimes called weights (“w”)
    - real-valued constants (can be positive or negative)
  - Input features  $x_1 \dots x_n$ ; define an additional constant input “1”
- A perceptron calculates 2 quantities:
  - 1. A weighted sum of the input features
  - 2. This sum is then thresholded by the  $T(\cdot)$  function
- Perceptron: a simple artificial model of human neurons
  - weights = “synapses”
  - threshold = “neuron firing”

# Perceptron Decision Boundary

- The perceptron is defined by the decision algorithm:

$$f(x; \theta) = \begin{cases} +1 & \text{if } \theta \cdot x^T > 0 \\ -1 & \text{otherwise} \end{cases}$$

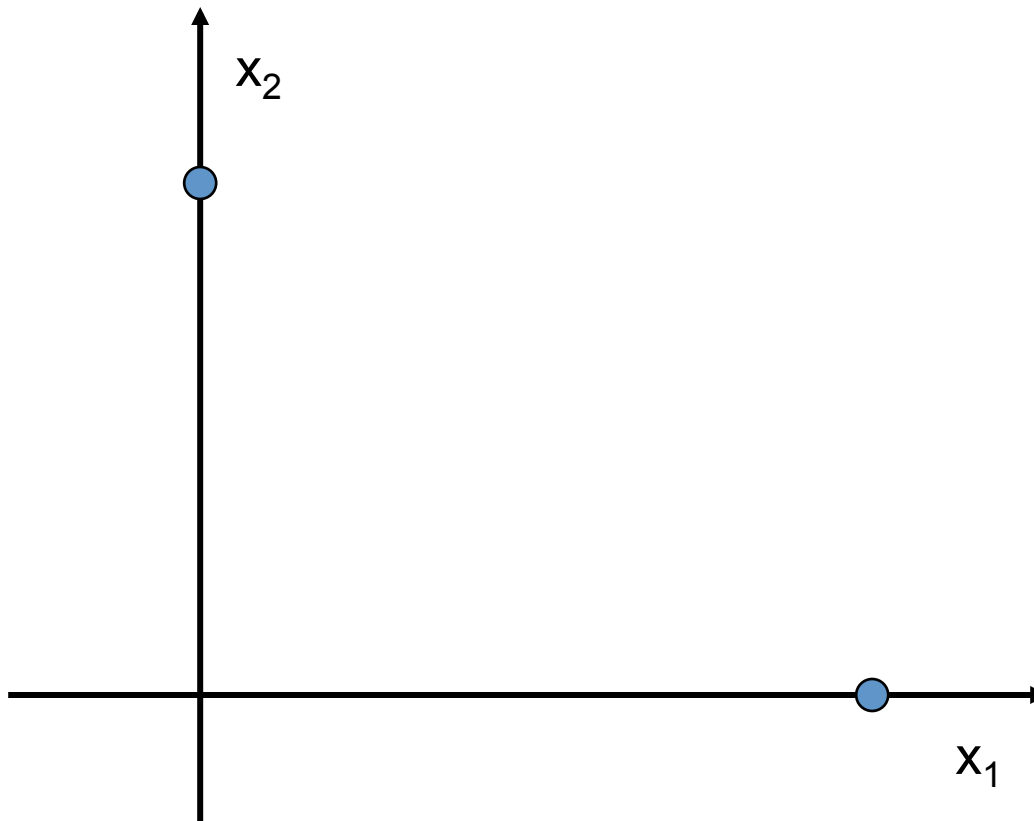
- The perceptron represents a hyperplane decision surface in d-dimensional space
  - A line in 2D, a plane in 3D, etc.
- The equation of the hyperplane is given by

$$\underline{\theta} \cdot \underline{\mathbf{x}}^T = 0$$

This defines the set of points that are on the boundary.

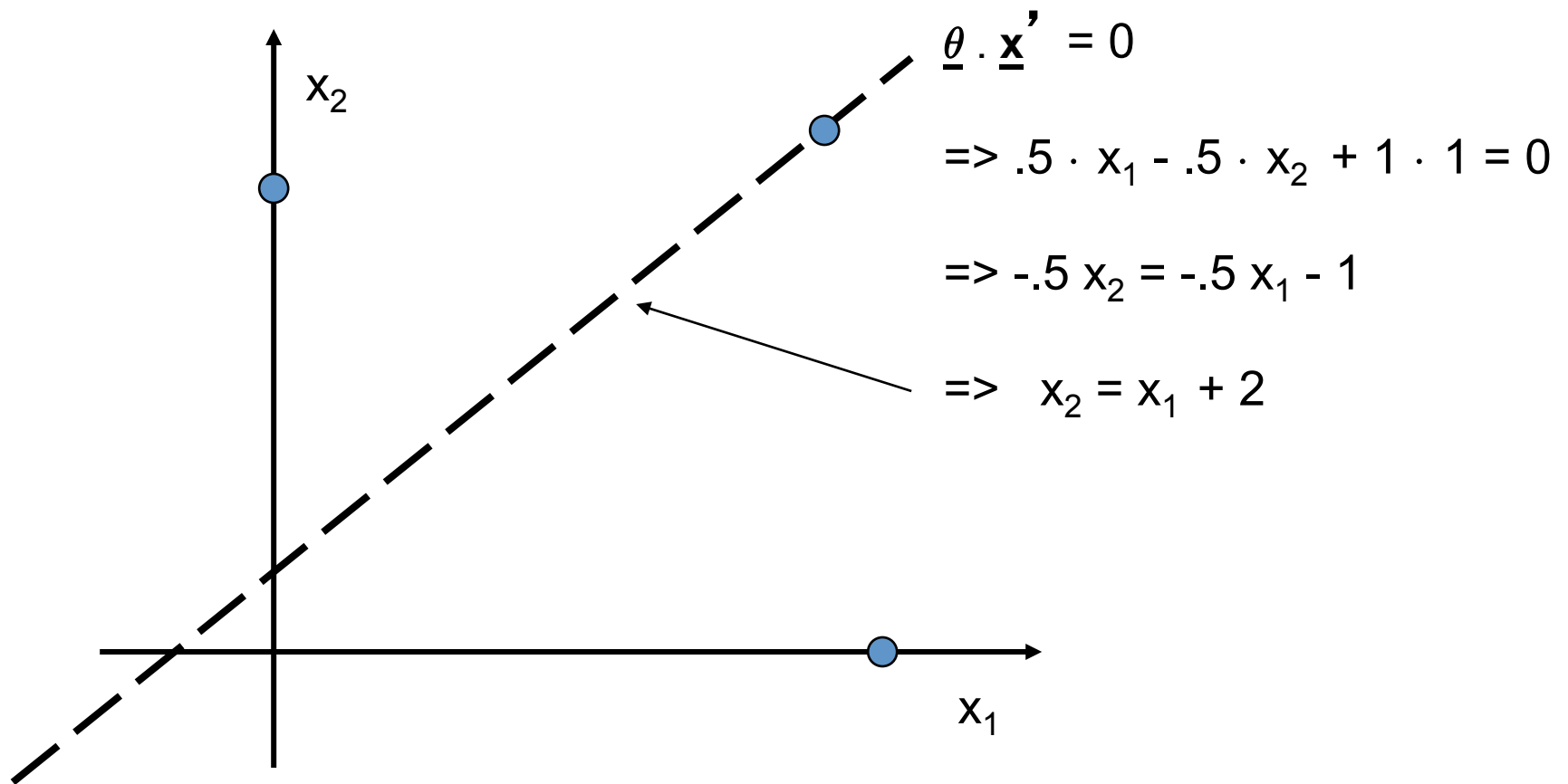
# Example, Linear Decision Boundary

$$\begin{aligned}\underline{\theta} &= (\theta_0, \theta_1, \theta_2) \\ &= (1, .5, -.5)\end{aligned}$$



# Example, Linear Decision Boundary

$$\begin{aligned}\underline{\theta} &= (\theta_0, \theta_1, \theta_2) \\ &= (1, .5, -.5)\end{aligned}$$



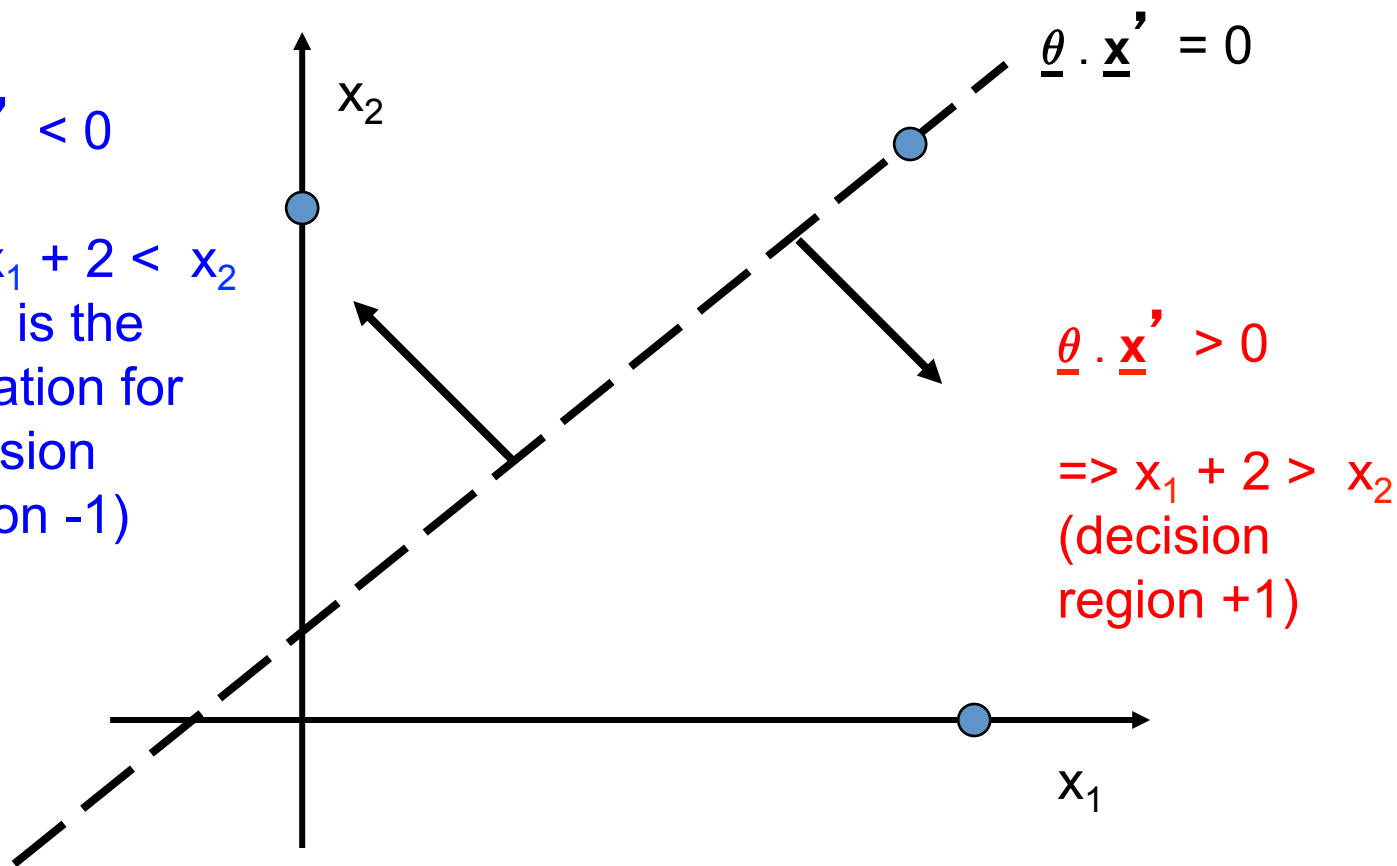


# Example, Linear Decision Boundary

$$\begin{aligned}\underline{\theta} &= (\theta_0, \theta_1, \theta_2) \\ &= (1, .5, -.5)\end{aligned}$$

$$\underline{\theta} \cdot \underline{x}' < 0$$

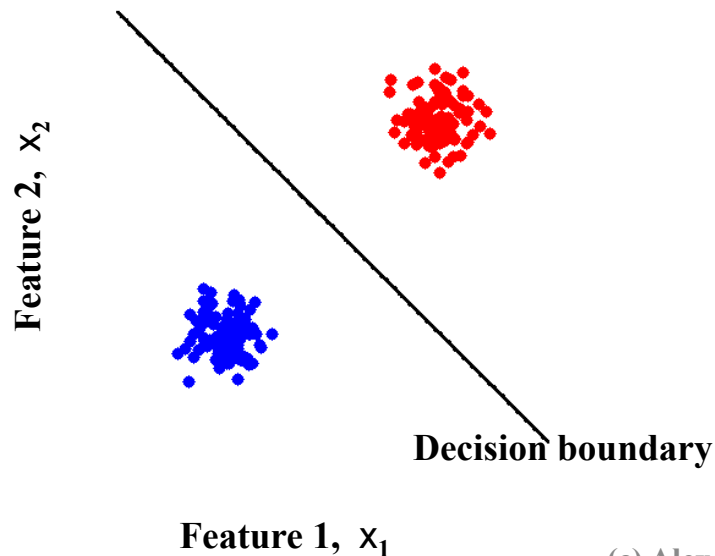
$\Rightarrow x_1 + 2 < x_2$   
(this is the  
equation for  
decision  
region -1)



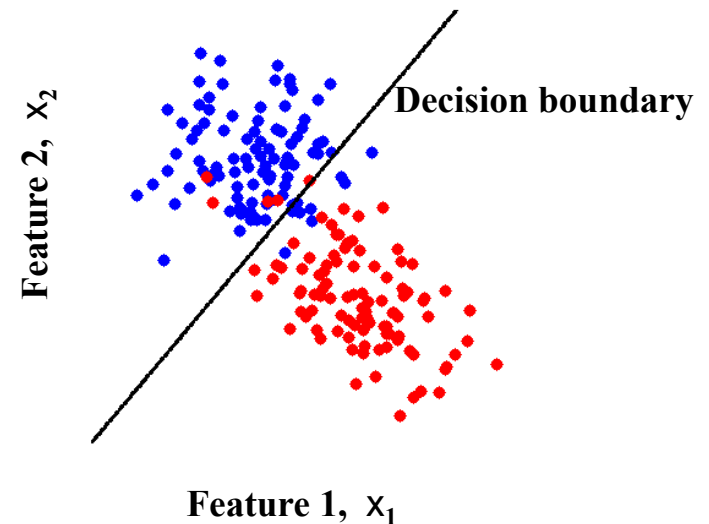
# Separability

- A data set is separable by a learner if
  - There is some instance of that learner that correctly predicts all the data points
- Linearly separable data
  - Can separate the two classes using a straight line in feature space
  - in 2 dimensions the decision boundary is a straight line

Linearly separable data

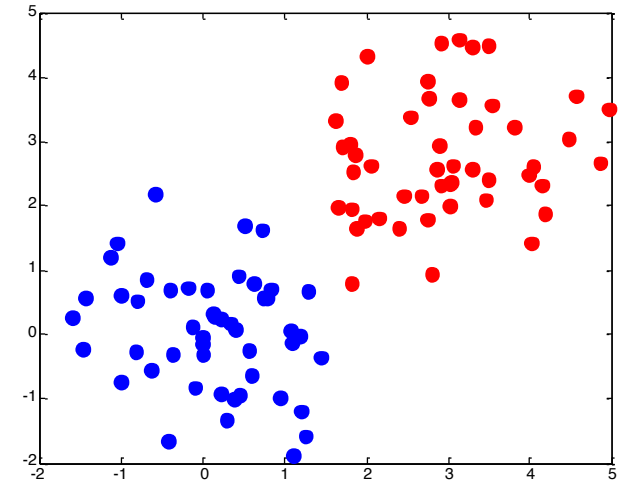


Linearly non-separable data

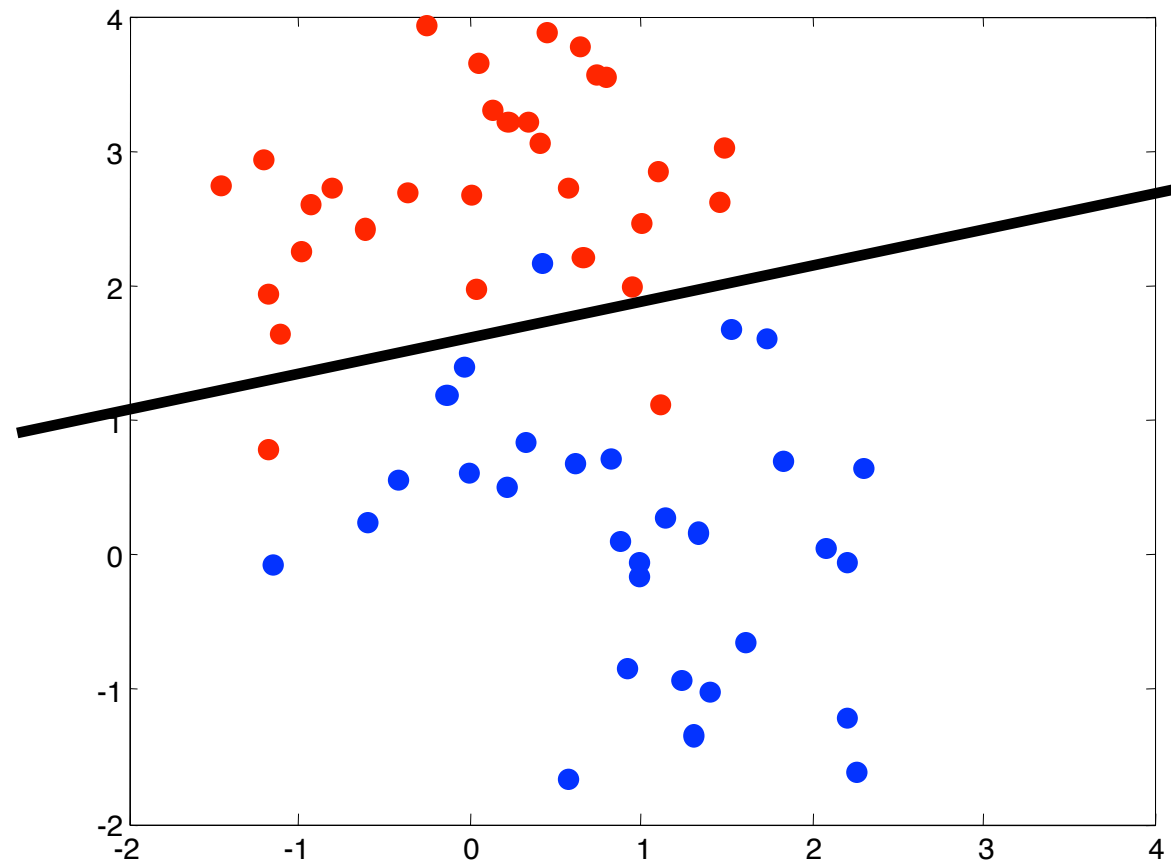


# Class overlap

- Classes may not be well-separated
- Same observation values possible under both classes
  - High vs low risk; features {age, income}
  - Benign/malignant cells look similar
  - ...
- Common in practice
- May not be able to perfectly distinguish between classes
  - Maybe with more features?
  - Maybe with more complex classifier?
- Otherwise, may have to accept some errors

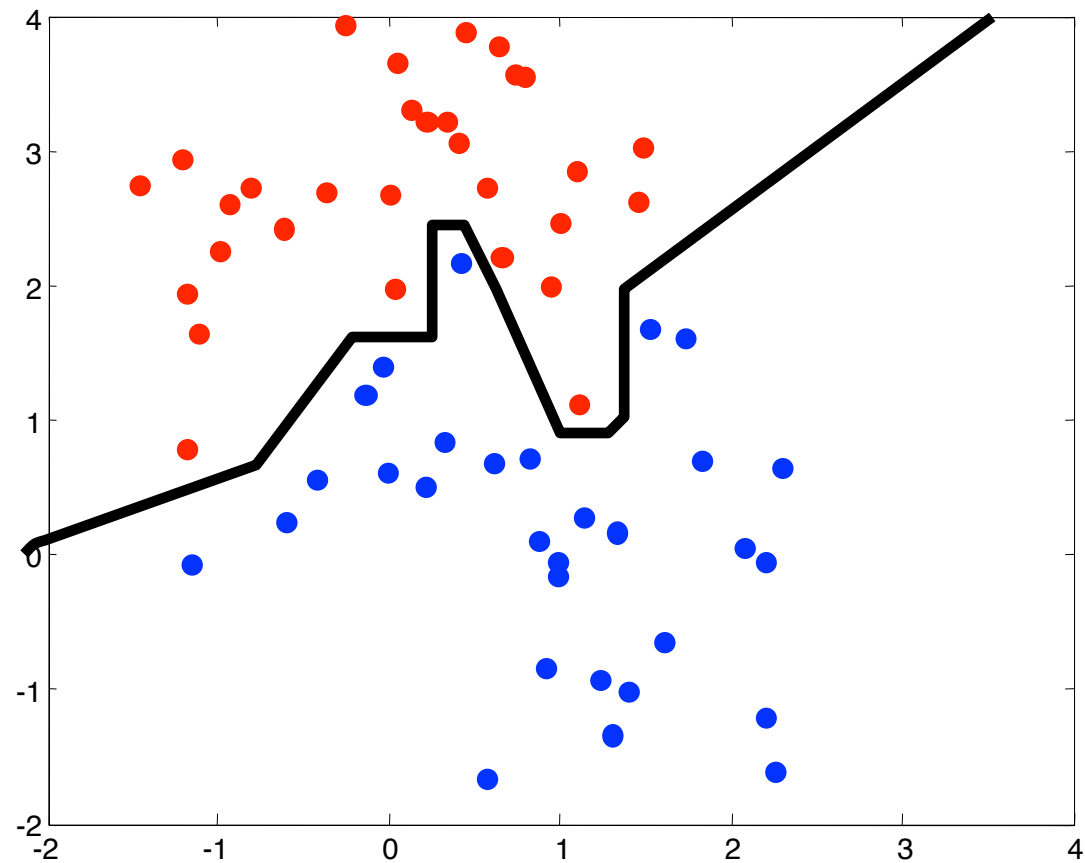


# Another example



(c) Alexander Ihler

# Non-linear decision boundary



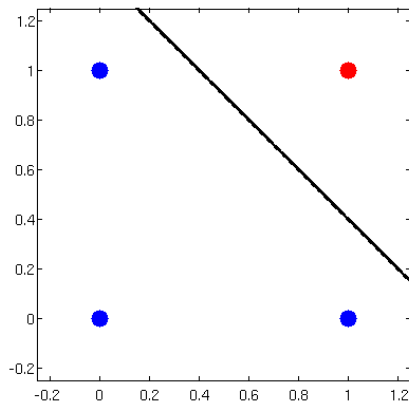
(c) Alexander Ihler

# Representational Power of Perceptrons

- What mappings can a perceptron represent perfectly?
  - A perceptron is a linear classifier
  - thus it can represent any mapping that is linearly separable
  - some Boolean functions like AND (on left)
  - but not Boolean functions like XOR (on right)

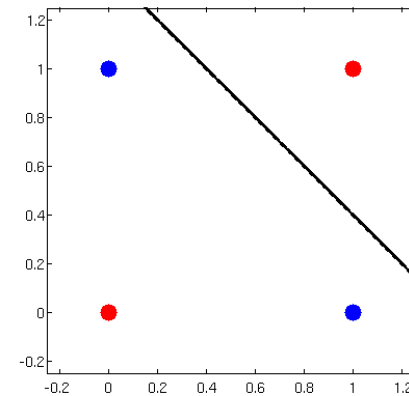
**“AND”**

$x_1$	$x_2$	$y$
0	0	-1
0	1	-1
1	0	-1
1	1	1



**“XOR”**

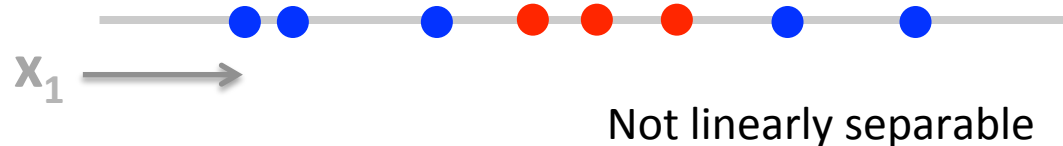
$x_1$	$x_2$	$y$
0	0	1
0	1	-1
1	0	-1
1	1	1



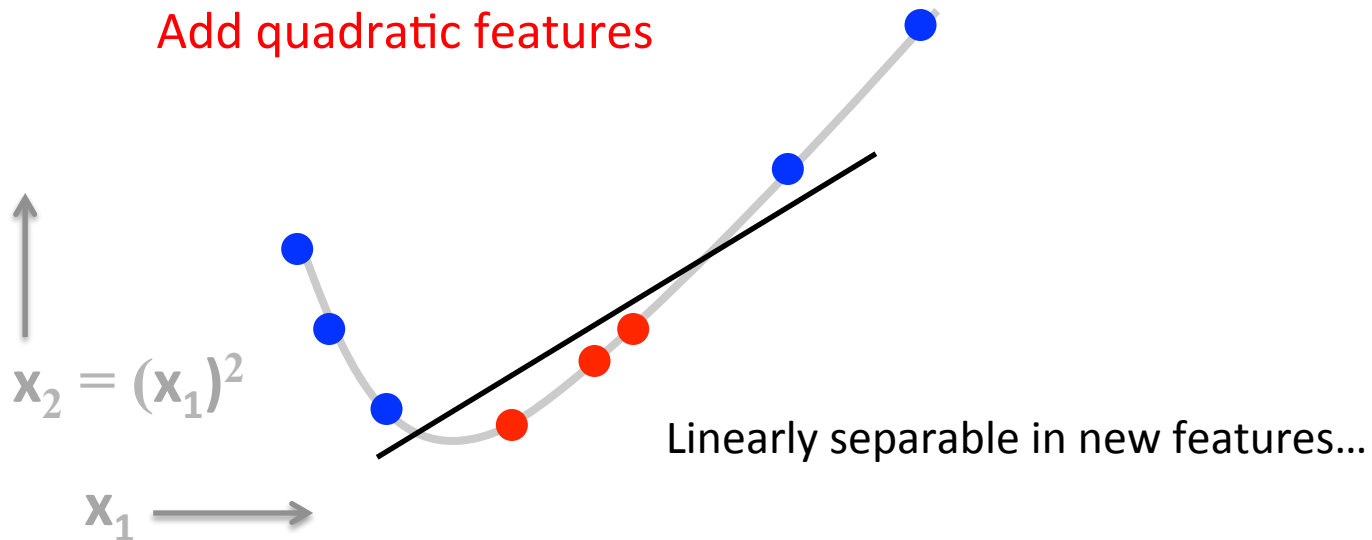
# Adding features

- Linear classifier can't learn some functions

1D example:



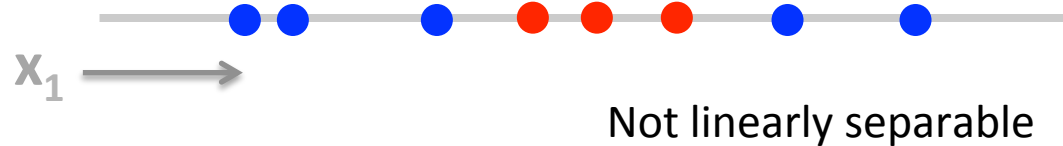
Add quadratic features



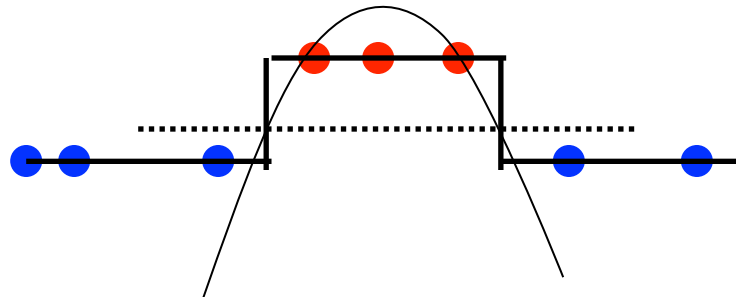
# Adding features

- Linear classifier can't learn some functions

1D example:



Quadratic features, visualized in original feature space:



$$y = T(a x^2 + b x + c)$$

More complex decision boundary:  $ax^2+bx+c = 0$

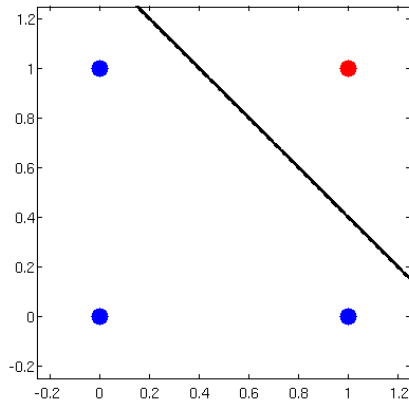


# Representational Power of Perceptrons

- What mappings can a perceptron represent perfectly?
  - A perceptron is a linear classifier
  - thus it can represent any mapping that is linearly separable
  - some Boolean functions like AND (on left)
  - but not Boolean functions like XOR (on right)

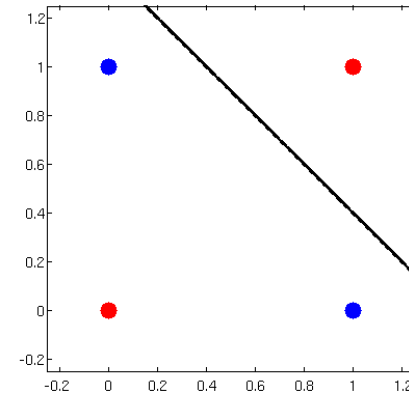
**“AND”**

$x_1$	$x_2$	$y$
0	0	-1
0	1	-1
1	0	-1
1	1	1



**“XOR”**

$x_1$	$x_2$	$y$
0	0	1
0	1	-1
1	0	-1
1	1	1



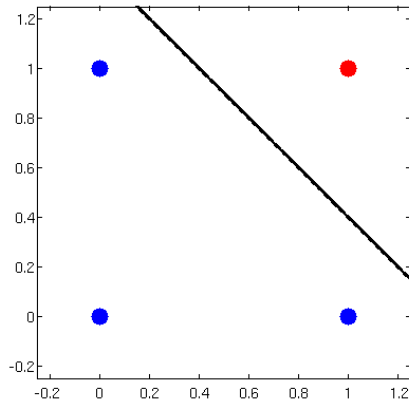
**What kinds of functions would we need to learn the data on the right?**

# Representational Power of Perceptrons

- What mappings can a perceptron represent perfectly?
  - A perceptron is a linear classifier
  - thus it can represent any mapping that is linearly separable
  - some Boolean functions like AND (on left)
  - but not Boolean functions like XOR (on right)

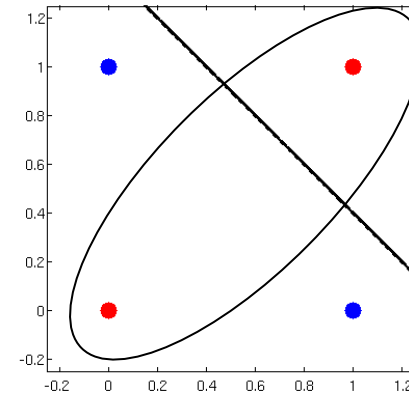
**“AND”**

$x_1$	$x_2$	$y$
0	0	-1
0	1	-1
1	0	-1
1	1	1



**“XOR”**

$x_1$	$x_2$	$y$
0	0	1
0	1	-1
1	0	-1
1	1	1



**What kinds of functions would we need to learn the data on the right?**

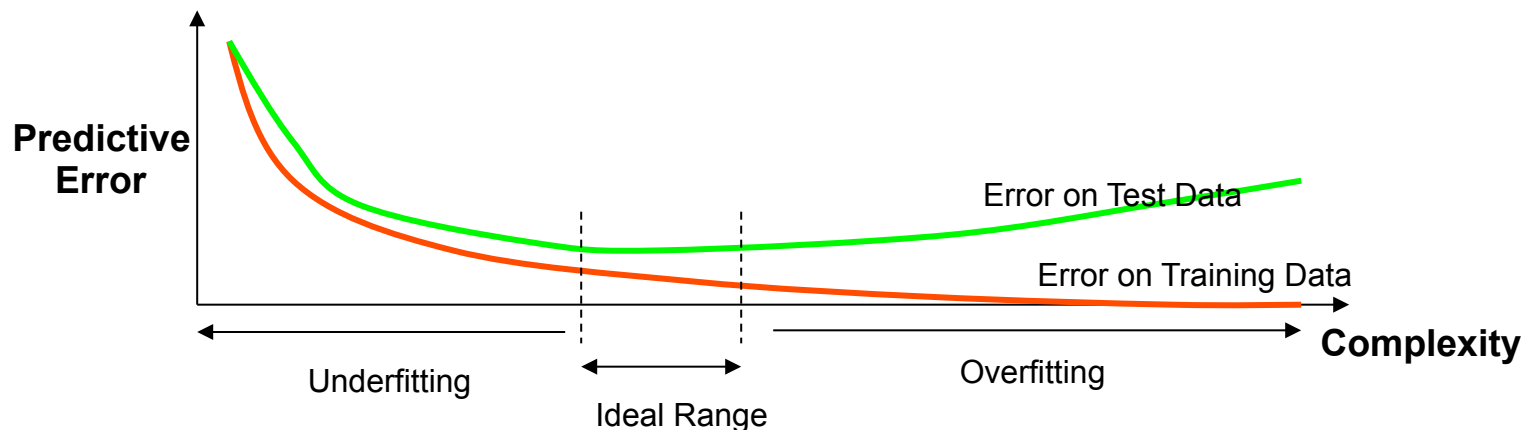
**Ellipsoidal decision boundary:**  $a x_1^2 + b x_1 + c x_2^2 + d x_2 + e x_1 x_2 + f = 0$

# Feature representations

- Features are used in a linear way
- Learner is dependent on representation
- Ex: discrete features
  - Mushroom surface: {fibrous, grooves, scaly, smooth}
  - Probably not useful to use  $x = \{1, 2, 3, 4\}$
  - Better: 1-of-K,  $x = \{ [1000], [0100], [0010], [0001] \}$
  - Introduces more parameters, but a more flexible relationship

# Effect of dimensionality

- Data are increasingly separable in high dimension – is this a good thing?
- “Good”
  - Separation is easier in higher dimensions (for fixed # of data  $m$ )
  - Increase the number of features, and even a linear classifier will eventually be able to separate all the training examples!
- “Bad”
  - Remember training vs. test error? Remember overfitting?
  - Increasingly complex decision boundaries can eventually get all the training data right, but it doesn't necessarily bode well for test data...



# Summary

---

- Linear classifier  $\Leftrightarrow$  perceptron
- Linear decision boundary
  - Computing and visualizing
- Separability
  - Limits of the representational power of a perceptron
- Adding features
  - Interpretations
  - Effect on separability
  - Potential for overfitting

+

# Machine Learning and Data Mining

## Linear classification: Learning

Prof. Alexander Ihler



# Learning the Classifier Parameters

- Learning from Training Data:
  - training data = labeled feature vectors
  - Find parameter values that predict well (low error)
    - error is estimated on the training data
    - “true” error will be on future test data
- Define a loss function  $J(\underline{\theta})$  :
  - Classifier error rate (for a given set of weights  $\underline{\theta}$  and labeled data)
- Minimize this loss function (or, maximize accuracy)
  - An optimization or search problem over the vector  $(\theta_1, \theta_2, \theta_0)$

# Training a linear classifier

- How should we measure error?
  - Natural measure = “fraction we get wrong” (error rate)

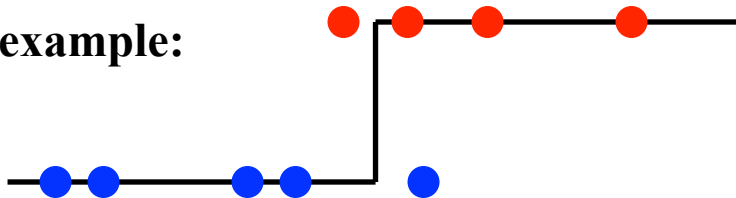
$$\text{err}(\theta) = \frac{1}{m} \sum_i \mathbb{1}[y^{(i)} \neq f(x^{(i)}; \theta)] \quad \text{where} \quad \mathbb{1}[y \neq \hat{y}] = \begin{cases} 1 & y \neq \hat{y} \\ 0 & \text{o.w.} \end{cases}$$

```
Yhat = np.sign( X.dot( theta.T ) );  
err = np.mean( Y != Yhat )
```

# predict class (+1/-1)  
# count errors: empirical error rate

- But, hard to train via gradient descent
  - Not continuous
  - As decision boundary moves, errors change abruptly

1D example:

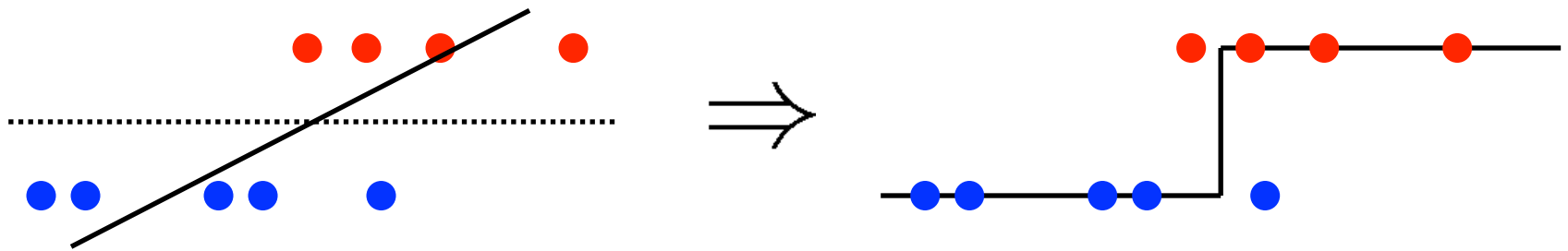


$$\begin{aligned} T(f) &= -1 & \text{if } f < 0 \\ T(f) &= +1 & \text{if } f > 0 \end{aligned}$$

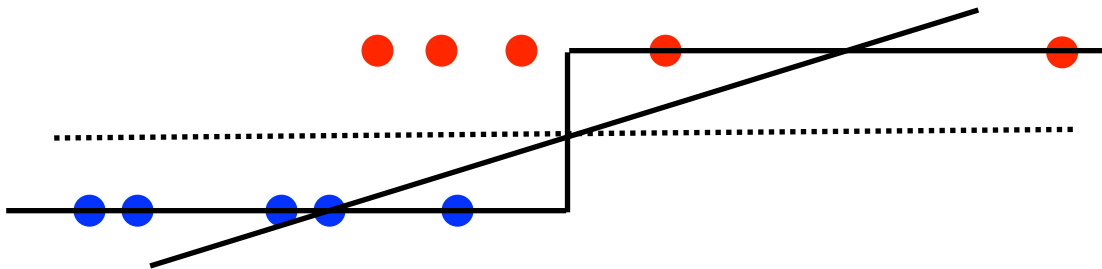


# Linear regression?

- Simple option: set  $\theta$  using linear regression



- In practice, this often doesn't work so well...
  - Consider adding a distant but “easy” point
  - MSE distorts the solution



# Perceptron algorithm

- Perceptron algorithm: an SGD-like algorithm

while  $\neg$  done:

for each data point  $j$ :

$$\hat{y}^{(j)} = \text{sign}(\theta \cdot x^{(j)}) \quad \text{(predict output for point } j\text{)}$$

$$\theta \leftarrow \theta + \alpha(y^{(j)} - \hat{y}^{(j)})x^{(j)} \quad \text{("gradient-like" step)}$$

- Compare to linear regression + MSE cost
  - Identical update to SGD for MSE except error uses thresholded  $\hat{y}(j)$  instead of linear response  $\underline{\theta} x'$  so:
    - (1) For correct predictions,  $y(j) - \hat{y}(j) = 0$
    - (2) For incorrect predictions,  $y(j) - \hat{y}(j) = \pm 2$

“adaptive” linear regression: correct predictions stop contributing

# Perceptron algorithm

- Perceptron algorithm: an SGD-like algorithm

while  $\neg$  done:

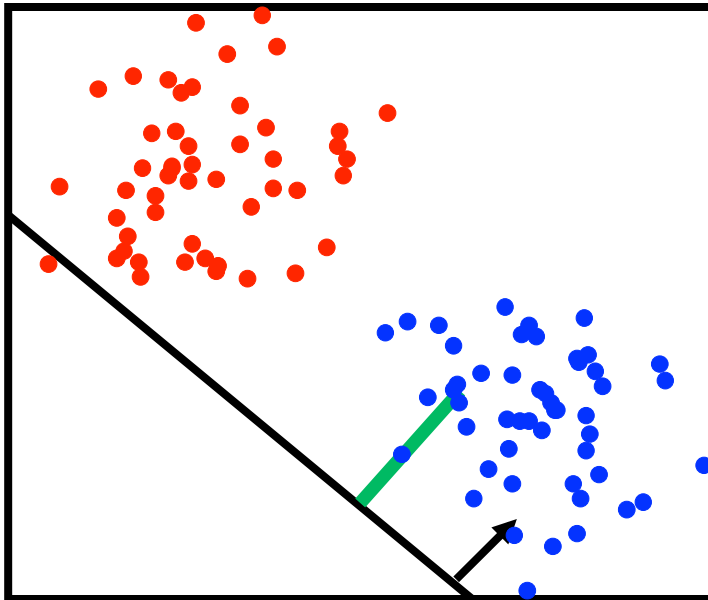
for each data point  $j$ :

$$\hat{y}^{(j)} = \text{sign}(\theta \cdot x^{(j)})$$

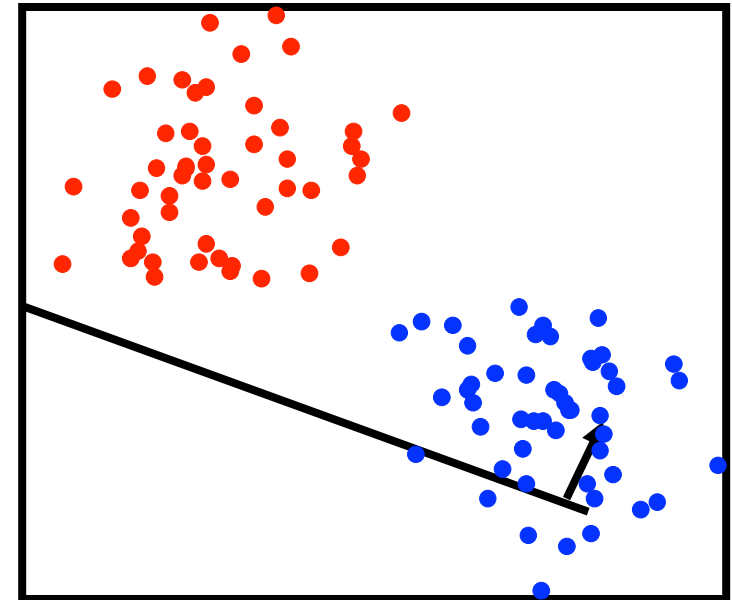
(predict output for point  $j$ )

$$\theta \leftarrow \theta + \alpha(y^{(j)} - \hat{y}^{(j)})x^{(j)}$$

("gradient-like" step)



$y(j)$   
predicted  
**incorrectly:**  
update  
weights



(c) Alexander Ihler

# Perceptron algorithm

- Perceptron algorithm: an SGD-like algorithm

while  $\neg$  done:

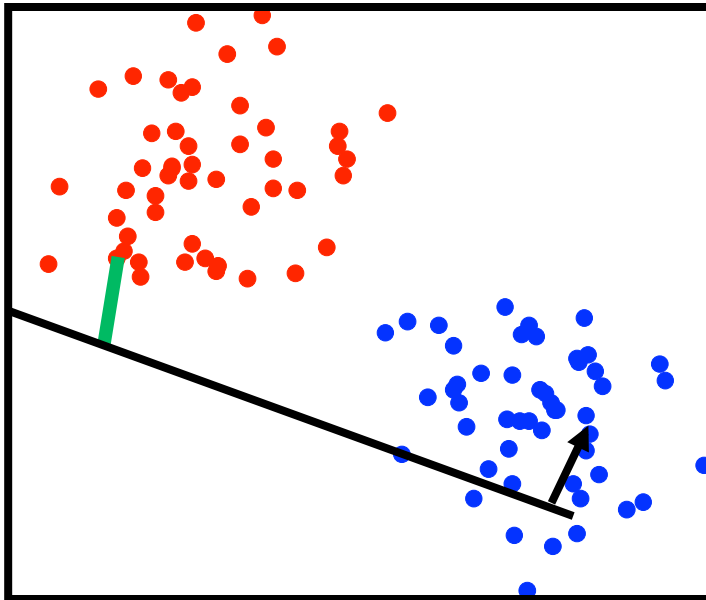
for each data point  $j$ :

$$\hat{y}^{(j)} = \text{sign}(\theta \cdot x^{(j)})$$

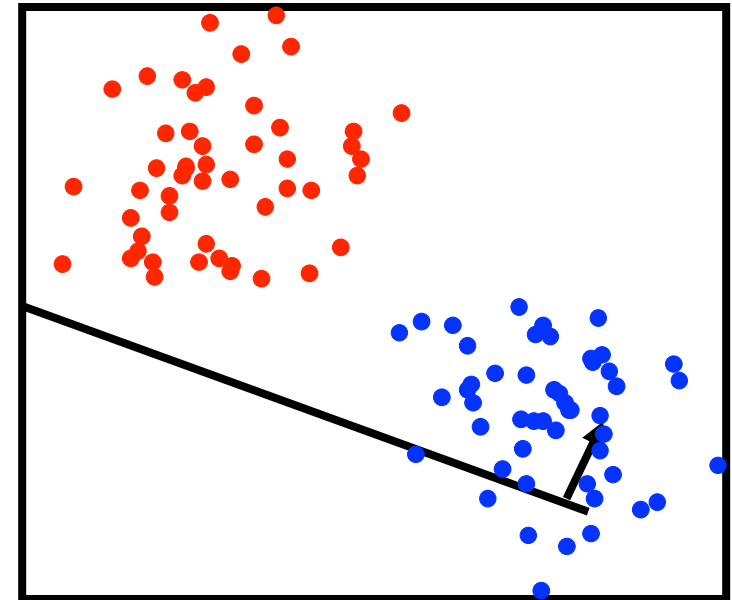
(predict output for point  $j$ )

$$\theta \leftarrow \theta + \alpha(y^{(j)} - \hat{y}^{(j)})x^{(j)}$$

(“gradient-like” step)



$y(j)$   
predicted  
**correctly:**  
no update



(c) Alexander Ihler

# Perceptron algorithm

- Perceptron algorithm: an SGD-like algorithm

while  $\neg$  done:

for each data point  $j$ :

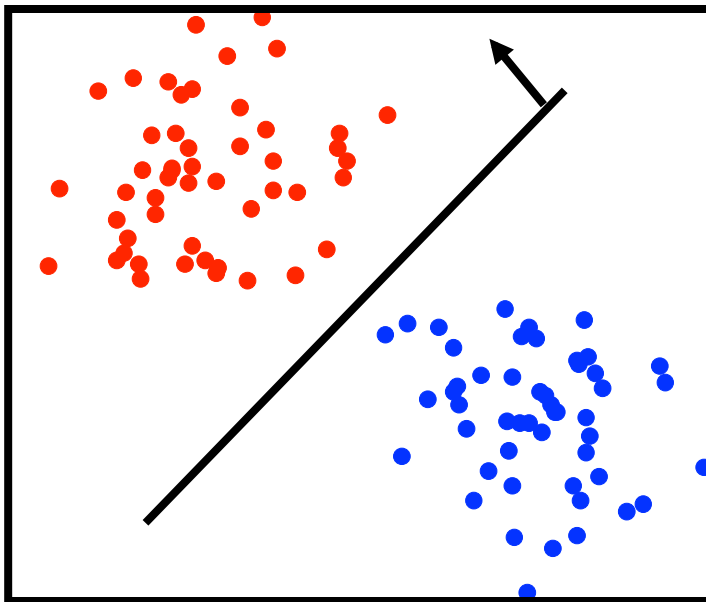
$$\hat{y}^{(j)} = \text{sign}(\theta \cdot x^{(j)})$$

(predict output for point  $j$ )

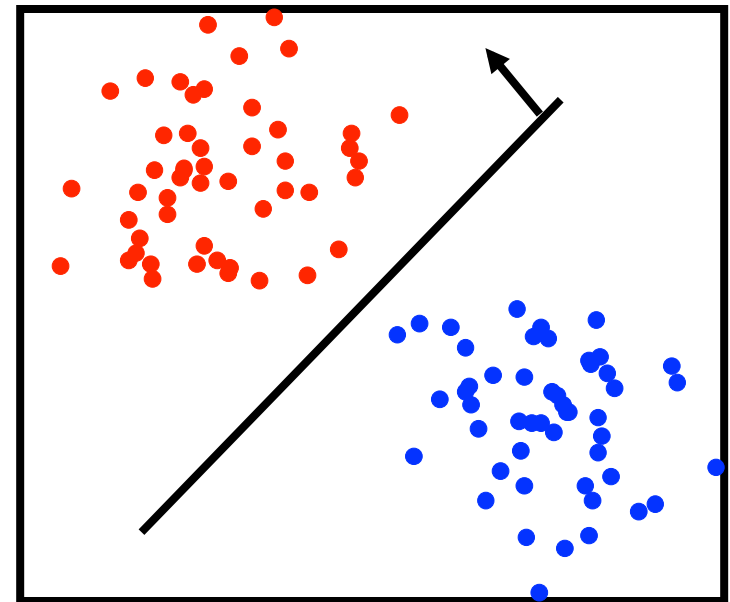
$$\theta \leftarrow \theta + \alpha(y^{(j)} - \hat{y}^{(j)})x^{(j)}$$

("gradient-like" step)

(Converges if data are linearly separable)



$y^{(j)}$   
predicted  
**correctly:**  
no update



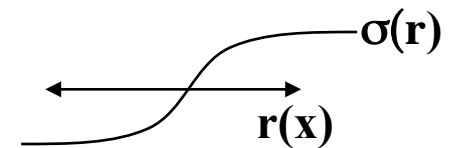
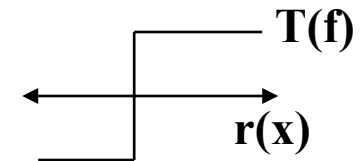
(c) Alexander Ihler

# Surrogate loss functions

- Another solution: use a “smooth” loss
  - e.g., approximate the threshold function
  - Usually some smooth function of distance
    - Example: logistic “sigmoid”, looks like an “S”
  - Now, measure e.g. MSE

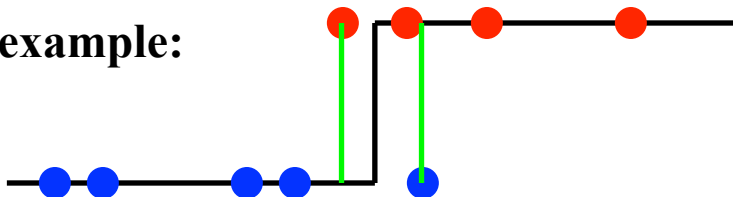
$$J(\underline{\theta}) = \frac{1}{m} \sum_j \left( \sigma(r(x^{(j)})) - y^{(j)} \right)^2$$

- Far from the decision boundary:  $|f(\cdot)|$  large, small error
- Nearby the boundary:  $|f(\cdot)|$  near 1/2, larger error

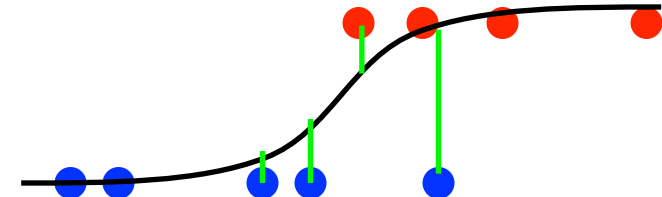


**Class  $y = \{0, 1\} \dots$**

1D example:



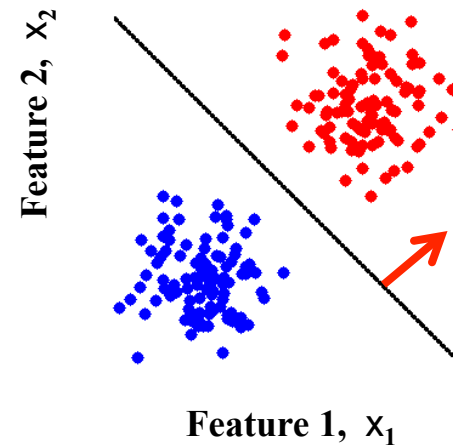
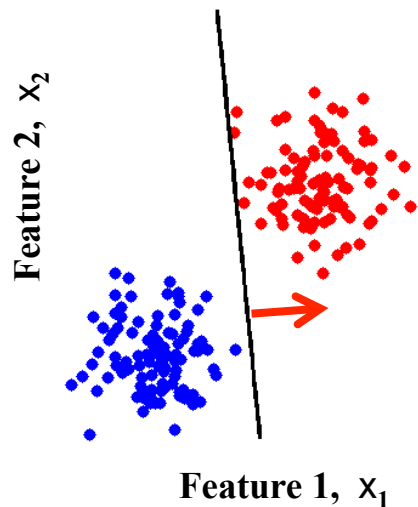
**Classification error = 2/9**



**MSE =  $(0^2 + 1^2 + .2^2 + .25^2 + .05^2 + \dots)/9$**

# Beyond misclassification rate

- Which decision boundary is “better”?
  - Both have zero training error (perfect training accuracy)
  - But, one of them seems intuitively better...



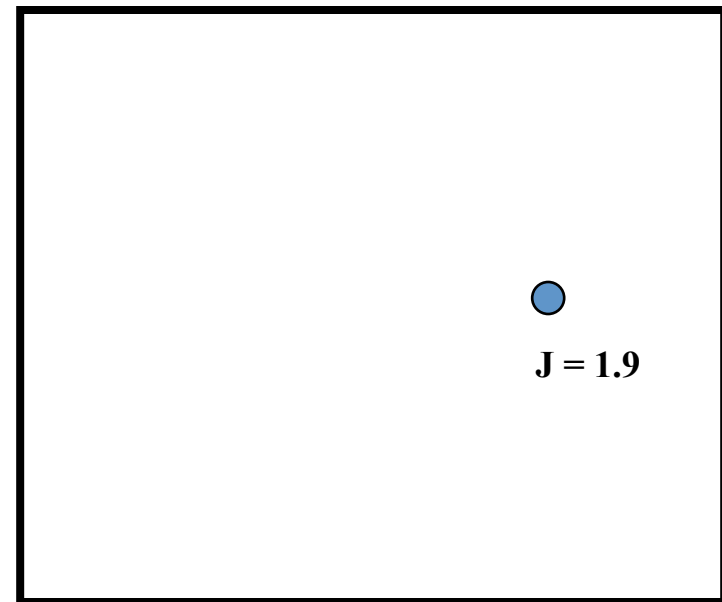
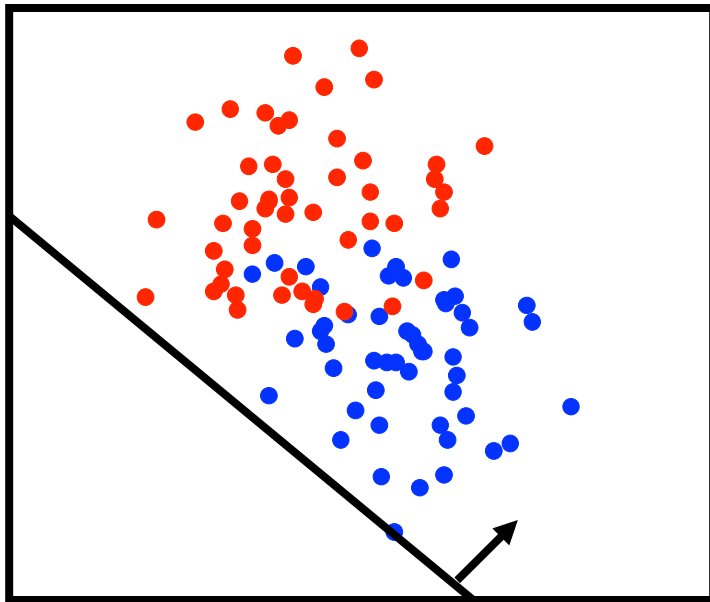
- Side benefit of many “smoothed” error functions
  - Encourages data to be far from the decision boundary
  - See more examples of this principle later...

# Training the Classifier

- Once we have a smooth measure of quality, we can find the “best” settings for the parameters of

$$r(x_1, x_2) = a \cdot x_1 + b \cdot x_2 + c$$

- Example: 2D feature space  $\Leftrightarrow$  parameter space



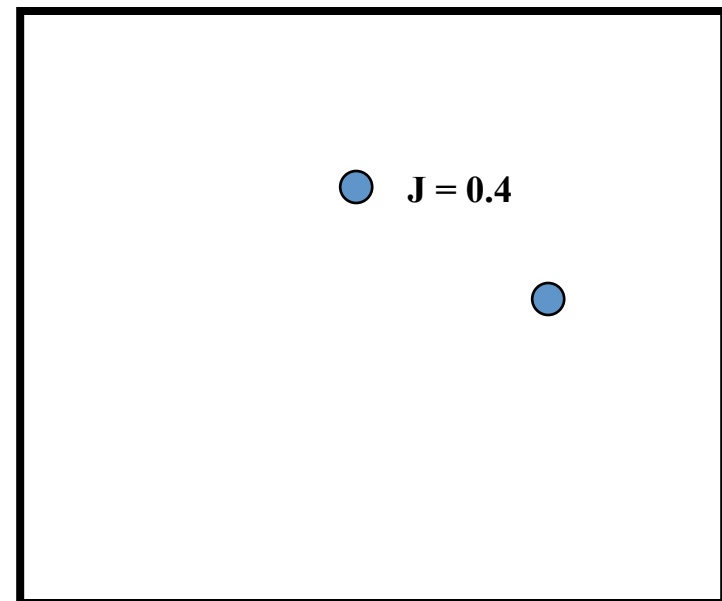
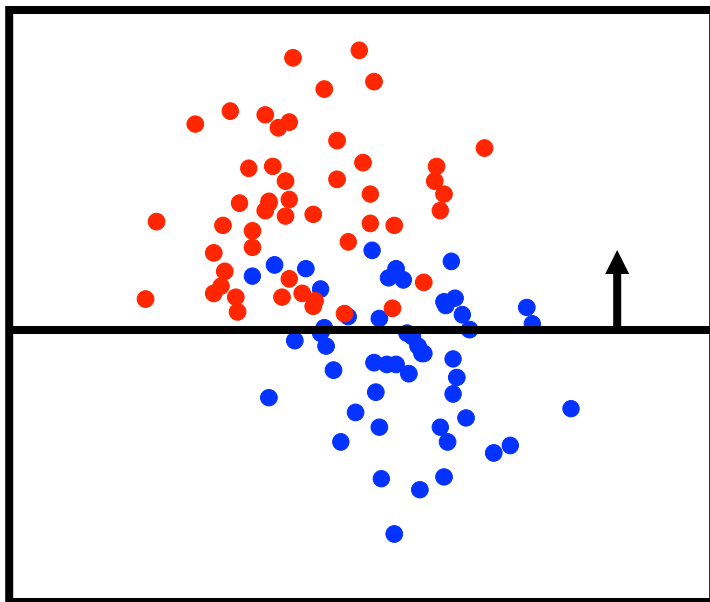


# Training the Classifier

- Once we have a smooth measure of quality, we can find the “best” settings for the parameters of

$$r(x_1, x_2) = a \cdot x_1 + b \cdot x_2 + c$$

- Example: 2D feature space  $\Leftrightarrow$  parameter space

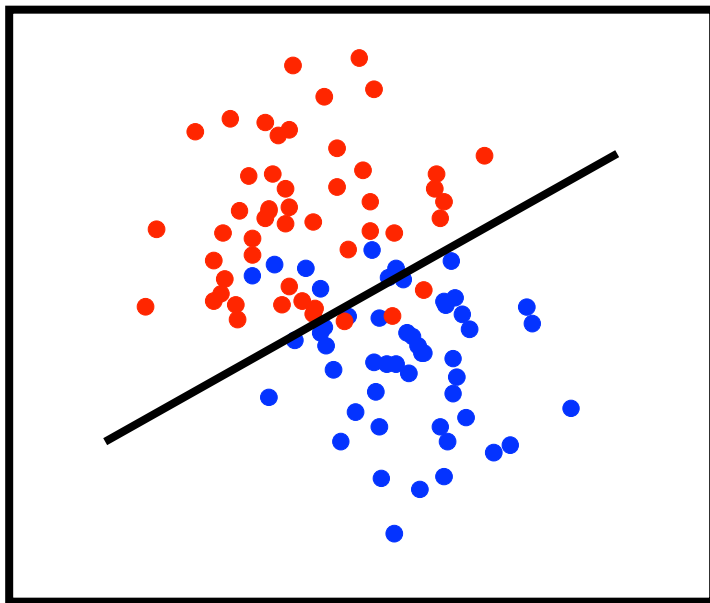


# Training the Classifier

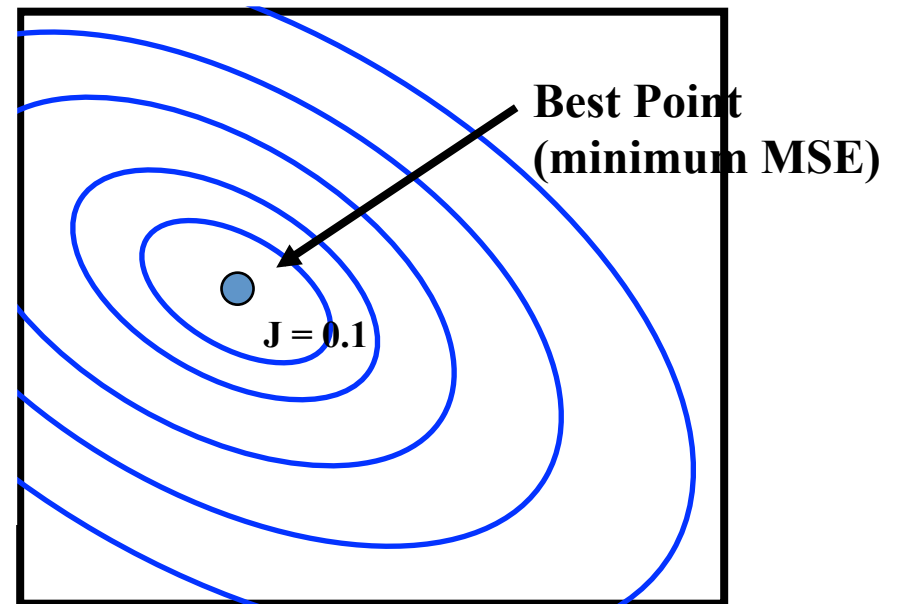
- Once we have a smooth measure of quality, we can find the “best” settings for the parameters of

$$r(x_1, x_2) = a \cdot x_1 + b \cdot x_2 + c$$

- Example: 2D feature space



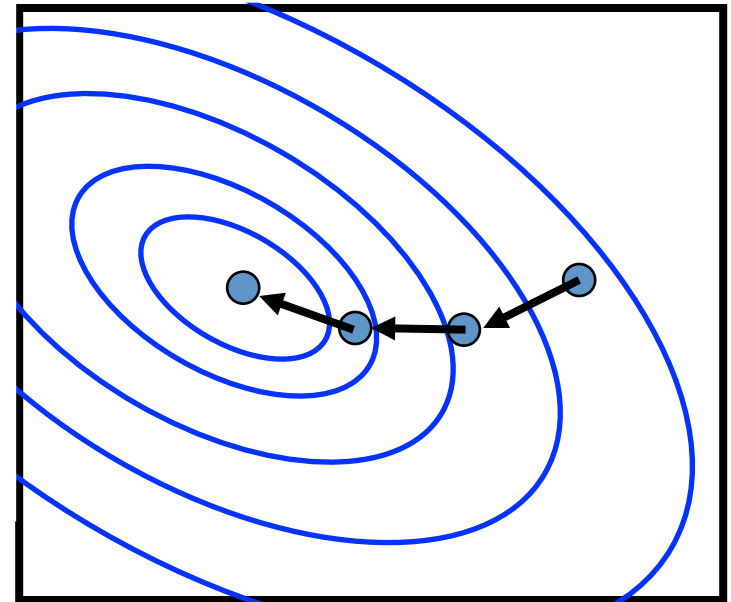
parameter space



# Finding the Best MSE

- As in linear regression, this is now just optimization
- Methods:
  - Gradient descent
    - Improve loss by small changes in parameters (“small” = learning rate)
  - Or, substitute your favorite optimization algorithm...
    - Coordinate descent
    - Stochastic search
    - Genetic algorithms

**Gradient Descent**



# Gradient Equations

- MSE (note, depends on function  $\sigma(\cdot)$ )

$$J(\underline{\theta} = [a, b, c]) = \frac{1}{m} \sum_i (\sigma(ax_1^{(i)} + bx_2^{(i)} + c) - y^{(i)})^2$$

- What's the derivative with respect to one of the parameters?
  - Recall the chain rule of calculus:

$$\frac{\partial}{\partial a} f(g(h(a))) = f'(g(h(a))) g'(h(a)) h'(a)$$

$$f(g) = (g)^2 \quad \Rightarrow \quad f'(g) = 2(g)$$

$$g(h) = \sigma(h) - y \quad \Rightarrow \quad g'(h) = \sigma'(h)$$

$$h(a) = ax_1^{(i)} + bx_2^{(i)} + c \quad \Rightarrow \quad h'(a) = x_1^{(i)}$$

w.r.t. b,c : similar;  
replace  $x_1$   
with  $x_2$  or 1

$$\frac{\partial J}{\partial a} = \frac{1}{m} \sum_i 2 \underbrace{(\sigma(\theta \cdot x^{(i)}) - y^{(i)})}_{\text{Error between class and prediction}} \underbrace{\partial \sigma(\theta \cdot x^{(i)}) x_1^{(i)}}_{\text{Sensitivity of prediction to changes in parameter "a"}}$$

# Saturating Functions

- Many possible “saturating” functions
- “Logistic” sigmoid (scaled for range [0,1]) is
$$\sigma(z) = 1 / (1 + \exp(-z))$$
- Derivative (slope of the function at a point  $z$ ) is
$$\partial\sigma(z) = \sigma(z) (1-\sigma(z))$$
- Python Implementation:

( $z$  = linear response,  $x^T\theta$ )

(to predict:  
threshold  $z$  at 0 or  
threshold  $\sigma(z)$  at  $\frac{1}{2}$ )

```
def sig(z):                                # logistic sigmoid
    return 1.0 / (1.0 + np.exp(-z)) # in [0,1]

def dsig(z):                               # its derivative at z
    return sig(z) * (1-sig(z))
```

For range  $[-1, +1]$ :

$$\rho(z) = 2 \sigma(z) - 1$$

$$\partial\rho(z) = 2 \sigma(z) (1-\sigma(z))$$

Predict: threshold  $z$  or  $\rho$  at zero

# Logistic regression

- Interpret  $\sigma(\underline{\theta} \cdot \mathbf{x}^T)$  as a probability that  $y = 1$
- Use a negative log-likelihood loss function
  - If  $y = 1$ , cost is  $-\log \Pr[y=1] = -\log \sigma(\underline{\theta} \cdot \mathbf{x}^T)$
  - If  $y = 0$ , cost is  $-\log \Pr[y=0] = -\log (1 - \sigma(\underline{\theta} \cdot \mathbf{x}^T))$
- Can write this succinctly:

$$J(\underline{\theta}) = -\frac{1}{m} \sum_i \underbrace{y^{(i)} \log \sigma(\underline{\theta} \cdot \mathbf{x}^{(i)})}_{\text{Nonzero only if } y=1} + \underbrace{(1-y^{(i)}) \log (1 - \sigma(\underline{\theta} \cdot \mathbf{x}^{(i)}))}_{\text{Nonzero only if } y=0}$$

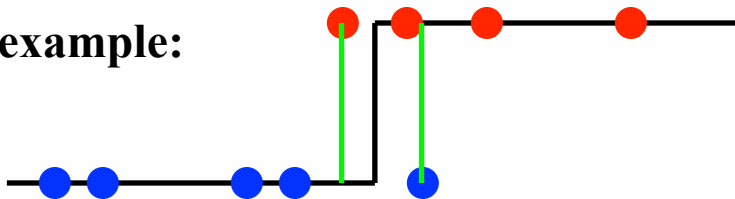
# Logistic regression

- Interpret  $\sigma(\underline{\theta} \cdot \mathbf{x}^T)$  as a probability that  $y = 1$
- Use a negative log-likelihood loss function
  - If  $y = 1$ , cost is  $-\log \Pr[y=1] = -\log \sigma(\underline{\theta} \cdot \mathbf{x}^T)$
  - If  $y = 0$ , cost is  $-\log \Pr[y=0] = -\log (1 - \sigma(\underline{\theta} \cdot \mathbf{x}^T))$
- Can write this succinctly:

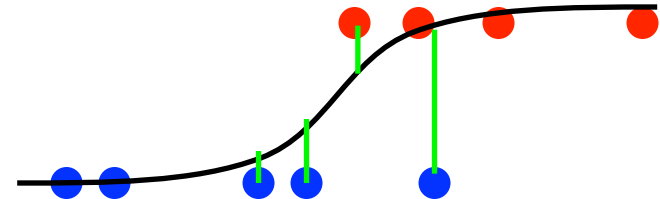
$$J(\underline{\theta}) = -\frac{1}{m} \sum_i y^{(i)} \log \sigma(\underline{\theta} \cdot \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\underline{\theta} \cdot \mathbf{x}^{(i)}))$$

- Convex! Otherwise similar: optimize  $J(\theta)$  via ...

1D example:



Classification error = MSE = 2/9



NLL =  $-(\log(.99) + \log(.97) + \dots)/9$

# Gradient Equations

- Logistic neg-log likelihood loss:

$$J(\underline{\theta}) = -\frac{1}{m} \sum_i y^{(i)} \log \sigma(\theta \cdot x^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(\theta \cdot x^{(i)}))$$

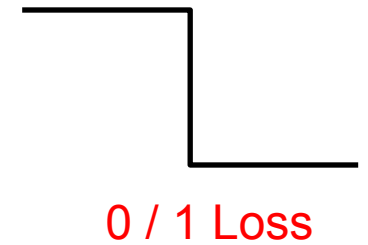
- What's the derivative with respect to one of the parameters?

$$\begin{aligned} \frac{\partial J}{\partial a} &= -\frac{1}{m} \sum_i y^{(i)} \frac{1}{\sigma(\theta \cdot x^{(i)})} \partial \sigma(\theta \cdot x^{(i)}) x_1^{(i)} + (1 - y^{(i)}) \dots \\ &= -\frac{1}{m} \sum_i y^{(i)} (1 - \sigma(\theta \cdot x^{(i)})) x_1^{(i)} - (1 - y^{(i)}) \dots \end{aligned}$$



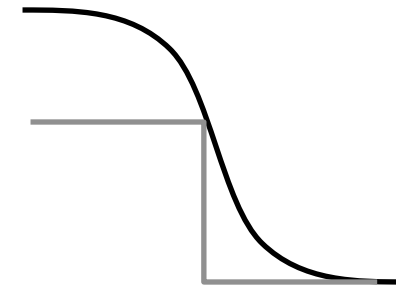
# Surrogate loss functions

- Replace 0/1 loss  $\Delta_i(\theta) = \mathbb{1}[T(\theta x^{(i)}) \neq y^{(i)}]$  with something easier:



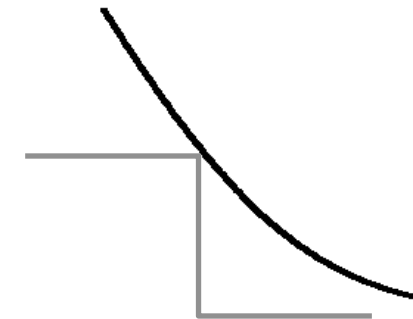
- Logistic MSE

$$J_i(\theta) = 4(\sigma(\theta x^{(i)}) - y^{(i)})^2$$



- Logistic Neg Log Likelihood

$$J_i(\underline{\theta}) = -\frac{y^{(i)}}{\log 2} \log \sigma(\theta \cdot x^{(i)}) + \dots$$



# Summary

- Linear classifier  $\Leftrightarrow$  perceptron
- Measuring quality of a decision boundary
  - Error rate (0/1 loss)
  - Logistic sigmoid + MSE criterion
  - Logistic Regression
- Learning the weights of a linear classifier from data
  - Reduces to an optimization problem
  - Perceptron algorithm
  - For MSE or Logistic NLL, we can do gradient descent
  - Gradient equations & update rules

# Multi-class linear models

- What about multiple classes? One option:

- Define one linear response per class
- Choose class with the largest response

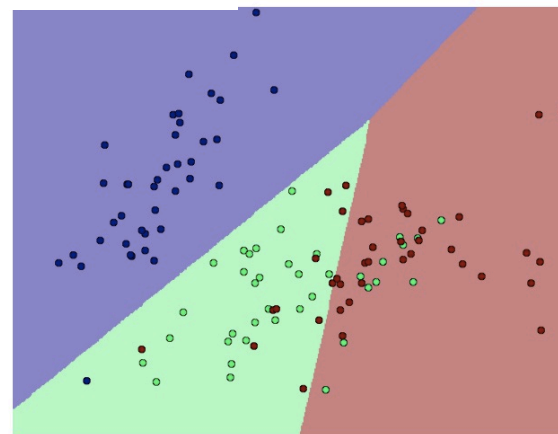
$$f(x; \theta) = \arg \max_c \theta_c \cdot x^T$$

$$\theta = \begin{bmatrix} \theta_{00} & \dots & \theta_{0n} \\ \vdots & \ddots & \vdots \\ \theta_{C0} & \dots & \theta_{Cn} \end{bmatrix}$$

- Boundary between two classes,  $c$  vs.  $c'$ ?

$$= \begin{cases} c & \text{if } \theta_c \cdot x^T > \theta_{c'} x^T \\ c' & \text{otherwise} \end{cases} \Leftrightarrow (\theta_c - \theta_{c'}) x^T > 0$$

- Linear boundary:  $(\theta_c - \theta_{c'}) x^T = 0$



# Multiclass linear models

- More generally, can define a generic linear classifier by

$$f(x; \theta) = \arg \max_y \theta \cdot \Phi(x, y)$$

- Example:  $y \in \{-1, +1\}$

$$\Phi(x, y) = y [1 \ x \ x^2 \ \dots]$$

$$f(x; \theta) = \begin{cases} +1 & \theta \cdot [1 \ x \ x^2 \ \dots] > -\theta \cdot [1 \ x \ x^2 \ \dots] \\ -1 & \text{o.w.} \end{cases}$$

(Standard perceptron rule)

# Multiclass linear models

- More generally, can define a generic linear classifier by

$$f(x; \theta) = \arg \max_y \theta \cdot \Phi(x, y)$$

- Example:  $y \in \{0, 1, 2, \dots\}$

$$\Phi(x, y) = [ \mathbb{1}[y = 0][1 \ x \ x^2 \ \dots] \ \mathbb{1}[y = 1][1 \ x \ x^2 \ \dots] \ \dots ]$$

$$\theta = [ \begin{bmatrix} \theta_{00} & \theta_{01} & \theta_{02} & \dots \end{bmatrix} \quad \begin{bmatrix} \theta_{10} & \theta_{11} & \theta_{12} & \dots \end{bmatrix} \quad \dots ]$$

(parameters for each class  $c$ )

$$f(x; \theta) = \arg \max_c \theta_c \cdot [1 \ x \ x^2 \ \dots]$$

(predict class with largest linear response)

# Multiclass perceptron algorithm

- Perceptron algorithm:
  - Make prediction  $f(x)$
  - Increase linear response of true target  $y$ ; decrease for prediction  $f$

While ( $\sim$ done)

For each data point  $j$ :

$$\begin{aligned} f^{(j)} &= \arg \max ( \underline{\theta}_c * \underline{x}^{(j)} ) && : \text{predict output for data point } j \\ \underline{\theta}_f &\leftarrow \underline{\theta}_f - \alpha \underline{x}^{(j)} && : \text{decrease response of class } f^{(j)} \text{ to } x^{(j)} \\ \underline{\theta}_y &\leftarrow \underline{\theta}_y + \alpha \underline{x}^{(j)} && : \text{increase response of true class } y^{(j)} \end{aligned}$$

— More general form update:

$$\begin{aligned} f(x; \theta) &= \arg \max_y \theta \cdot \Phi(x, y) \\ \theta &\leftarrow \theta + \alpha (\Phi(x, y) - \Phi(x, f(x))) \end{aligned}$$

# Multilogit regression

- Define the probability of each class:

$$p(Y = y|X = x) = \frac{\exp(\theta_y \cdot x^T)}{\sum_c \exp(\theta_c \cdot x^T)}$$

(Y binary = logistic regression)

- Then, the NLL loss function is:

$$J(\theta) = -\frac{1}{m} \sum_i \log p(y^{(i)}|x^{(i)}) = -\frac{1}{m} \sum_i \left[ \theta_{y^{(i)}} \cdot x^{(i)} - \log \sum_c \exp(\theta_c \cdot x^{(i)}) \right]$$

- P: “confidence” of each class
  - Soft decision value
- Decision: predict most probable
  - Linear decision boundary
- Convex loss function

