

# Machine Learning

## **Multiclass classification and Learning as Optimization**

Dan Goldwasser

[dgoldwas@purdue.edu](mailto:dgoldwas@purdue.edu)

# Multi-Categorical Output Tasks

- So far, our discussion was limited to ***binary predictions***
  - Well, ***almost*** (?)
- What happens if our decision is not over binary labels?
  - ***Many interesting classification problems are not!***
  - **Credit card**: Approved, Deny, Further investigation needed
  - **Document classification**: sports, finance, politics
  - **OCR**: 0,1,2,3..9,A,..,Z

***How can we approach these problems?***

- ***We will look into different reductions to binary classification problems!***

***Hint:*** What is the computer science solution to: “I can solve problem A, but now I have problem B, so...”

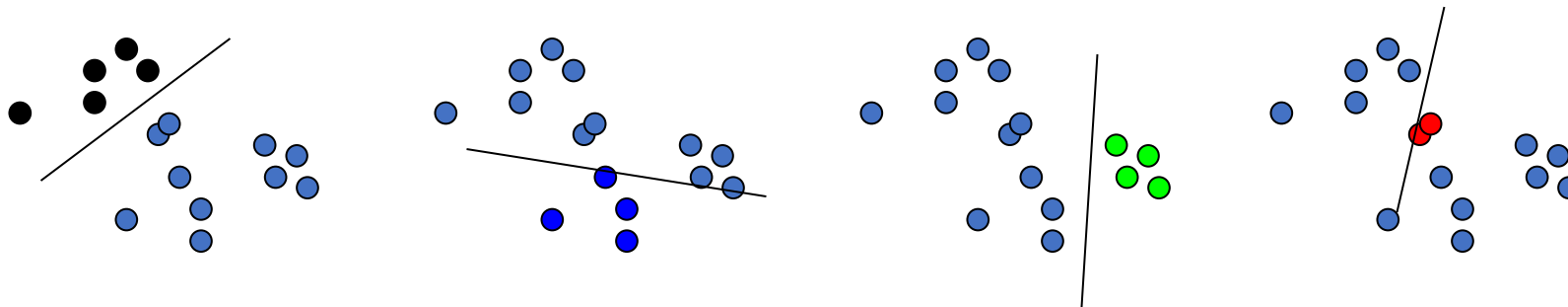
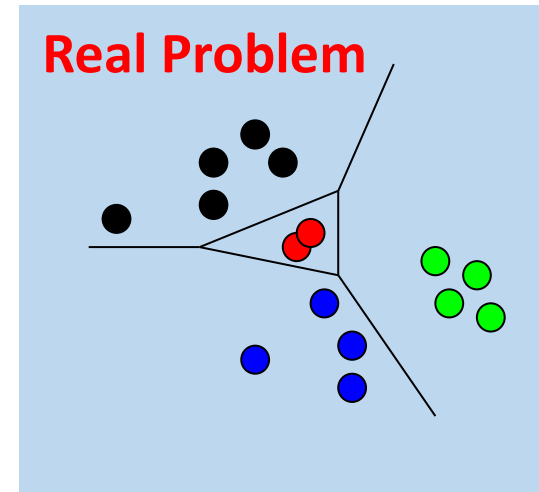
# Learning via One-Versus-All

- Find  $v_r, v_b, v_g, v_y \in \mathbf{R}^n$  such that

- $v_r x > 0$  iff  $y = \text{red}$   $\otimes$
- $v_b x > 0$  iff  $y = \text{blue}$   $\checkmark$
- $v_g x > 0$  iff  $y = \text{green}$   $\checkmark$
- $v_y x > 0$  iff  $y = \text{yellow}$   $\checkmark$

- Classification:*  $f(x) = \operatorname{argmax}_i (v_i x)$

$$\mathbf{H} = \mathbf{R}^{kn}$$

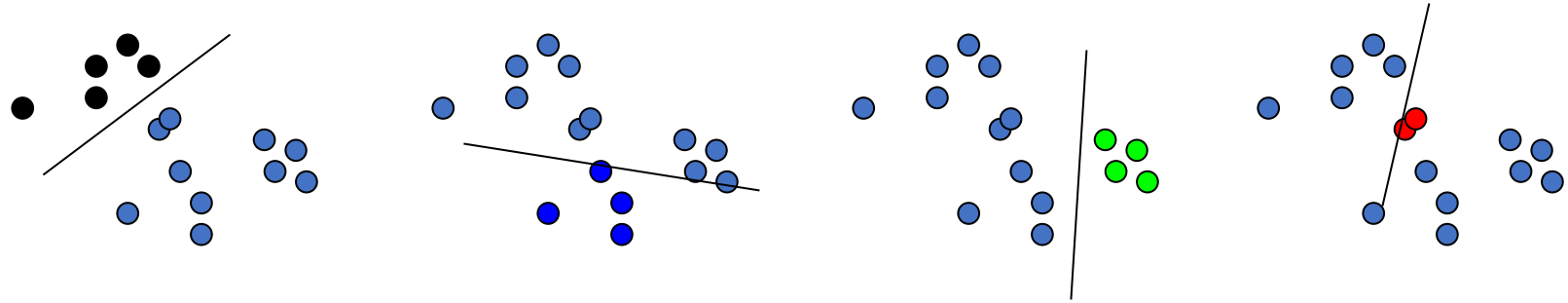


# Problems with Decompositions

- Learning optimizes over *local* metrics
  - Poor *global* performance
  - What is the metric?
    - We don't care about *local* classifiers performance
  - ***Poor decomposition  $\Rightarrow$  poor performance***
  - Especially true for Error Correcting Output Codes
- **Difficulty:** *how to ensure that the resulting problems are separable.*
- Can we ensure separability and learn the real objective?
- Real objective: final performance, rather than local metric

# Intuition

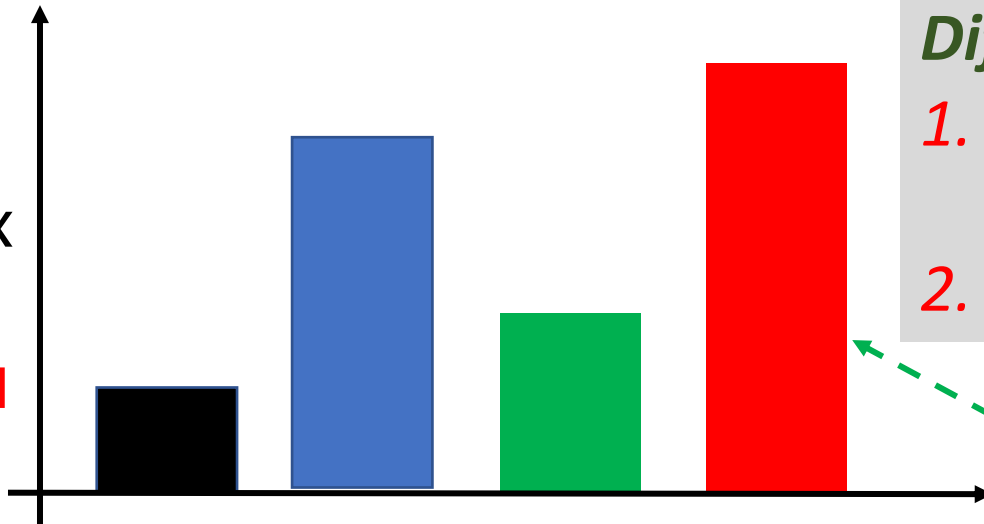
- One vs. All



- *We can think about the problem in a different way*

**Score =**  
 $F_L(x) = W_L x$

Each  $W_L$  is defined  
over different  
features



*Difference compared to 1vs.All*

1. Define the learning problem for **each label** over different feature space.
2. "Global" training objective

*Highest score =  
predicted label*

# Linear Separability for Multiclass

- We are learning  $k$   $n$ -dimensional weight vectors, so we can concatenate the  $k$  weight vectors into  $\mathbf{w} = (w_1, w_2, \dots, w_k) \in \mathbb{R}^{nk}$
- **Key Construction:** (Kesler Construction)
- Represent each example  $(\mathbf{x}, y)$ , as an  $nk$ -dimensional vector,  $\mathbf{x}_y$  with  $\mathbf{x}$  embedded in the  $y$ -th part of it ( $y=1, 2, \dots, k$ ) and the other coordinates are 0

E.g.,  $\mathbf{x}_y = (\mathbf{0}, \mathbf{x}, \mathbf{0}, \mathbf{0}) \in \mathbb{R}^{kn}$  (here  $k=4, y=2$ )

**Example:** predict sentiment of a product review

*Features – words (unigram)*

*Labels: Good, Bad, Neutral.*

The feature corresponding  
to “sick” when the label is “good”

( [“sick”,good], ..., [“sick”,neutral], ..., [“sick”,bad],...)

Good label features

Neutral label features

bad label features

# Linear Separability for Multiclass

- We are learning  $k$   $n$ -dimensional weight vectors, so we can concatenate the  $k$  weight vectors into  $\mathbf{w} = (w_1, w_2, \dots, w_k) \in \mathbb{R}^{nk}$
- **Key Construction:** (Kesler Construction)
- Represent each example  $(\mathbf{x}, y)$ , as an  $nk$ -dimensional vector,  $\mathbf{x}_y$  with  $\mathbf{x}$  embedded in the  $y$ -th part of it ( $y=1, 2, \dots, k$ ) and the other coordinates are 0

E.g.,  $\mathbf{x}_y = (0, \mathbf{x}, 0, 0) \in \mathbb{R}^{kn}$  (here  $k=4, y=2$ )

- Now we can understand the decision
  - Predict  $y$  iff  $\forall y' \in \{1, \dots, k\}, y \neq y' \quad (\mathbf{w}_y^T - \mathbf{w}_{y'}^T) \mathbf{x} > 0$

**Conclusion:** The set  $(\mathbf{x}_{yy'}, +) \equiv (\mathbf{x}_y - \mathbf{x}_{y'}, +)$  is *linearly separable* from the set  $(\mathbf{x}_{yy'}, -)$  using the linear separator  $\mathbf{w} \in \mathbb{R}^{kn}$ ,

# Learning via Kesler Construction

- A **perceptron update** rule applied in the **nk-dimensional space** due to a mistake in  $w^T x_{ij} > 0$
- Implies the following update:

- Given example  $(x, i)$  (example  $x \in \mathbb{R}^n$ , labeled  $i$ )
  - $\forall j = 1, \dots, k, i \neq j$
  - If  $(w_i^T - w_j^T) x < 0$  (mistaken prediction)
    - $w_i \leftarrow w_i + x$  (promotion)
    - $w_j \leftarrow w_j - x$  (demotion)

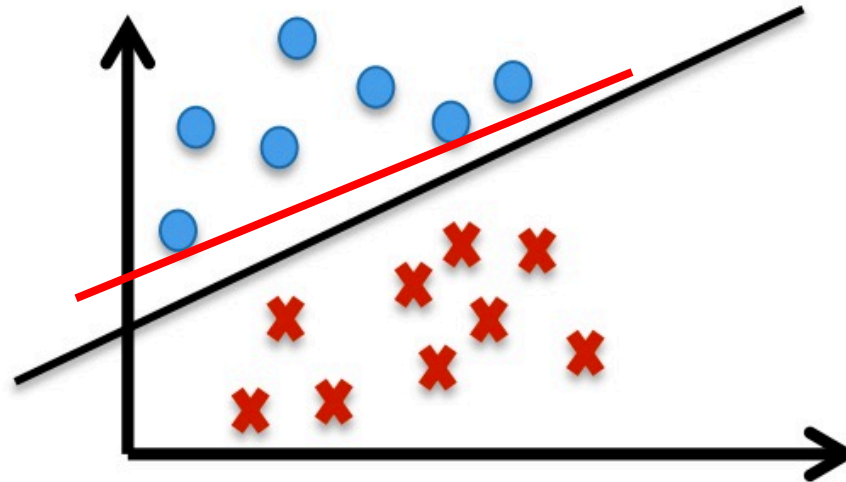
**Conservative version:** only demote the the weights corresponding to the label that go the highest score (incorrectly)

For any given example, you can potentially make  $K-1$  promotion steps for  $w_i$  and  $K-1$  demotion steps, for the *different*  $w_j$



# Learning using the Perceptron Algorithm

- Perceptron guarantee: *find a linear separator (if the data is separable)*



- There could be many models consistent with the training data
  - ***How did the perceptron algorithm deal with this problem?***
- Trading some training error for better generalization

# Reminder: Loss functions

- To formalize performance let's define a *loss function*:

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

- Where  $\hat{y}$  is the gold label
- *The loss function measures the error on a single instance*
  - Specific definition depends on the learning task

## ***Regression***

$$loss(y, \hat{y}) = (y - \hat{y})^2$$

## ***Binary classification***

$$loss(y, \hat{y}) = \begin{cases} 0 & y = \hat{y} \\ 1 & \text{otherwise} \end{cases}$$

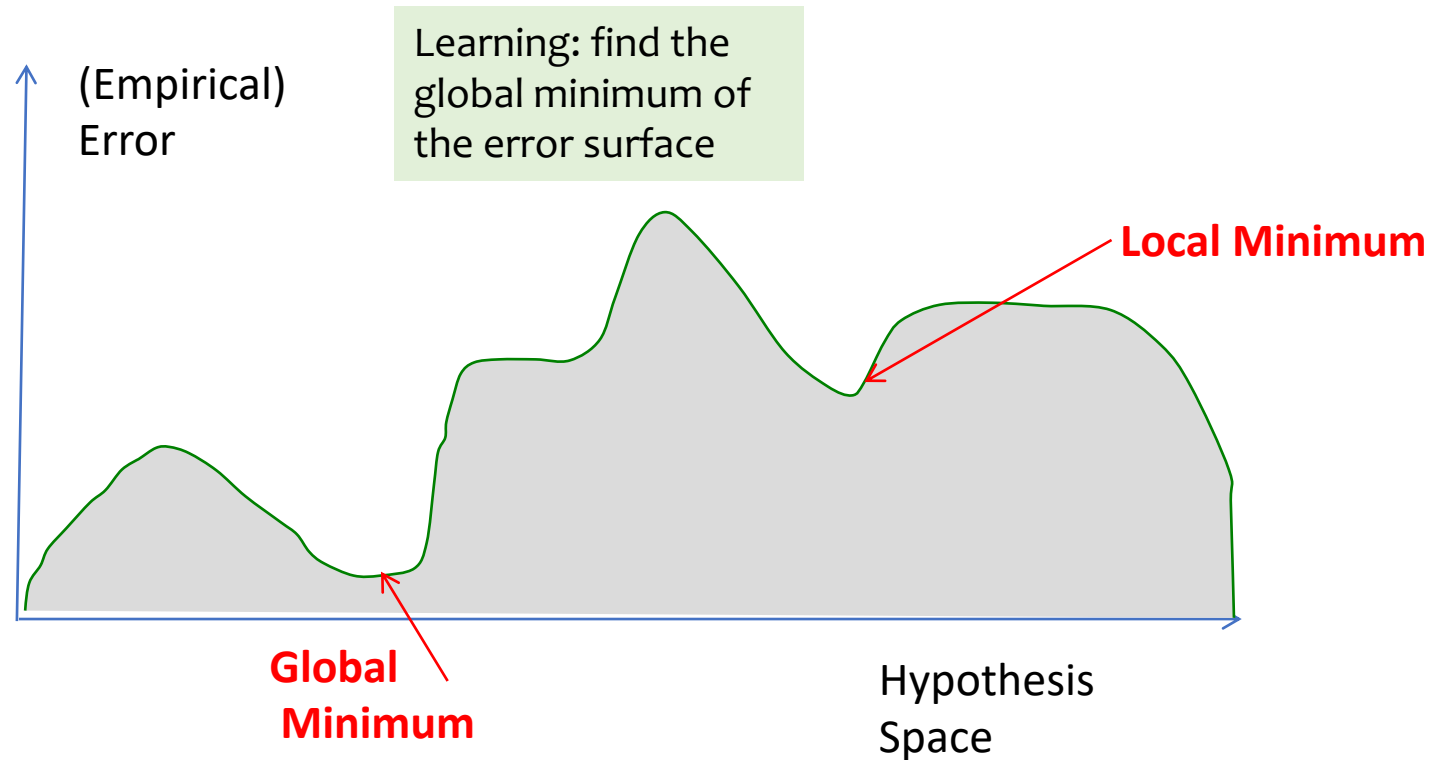
# Empirical Risk Minimization

- Learning: Given a training dataset  $D$ , return the classifier  $f(x)$  that minimizes the empirical risk
- *Given this definition we can view learning as a **minimization problem***
  - The objective function to minimize (empirical risk) is defined with respect to a specific loss function
  - Our minimization procedure (aka **learning**) will be influenced by the choice of loss function
    - **Some are easier to minimize than other!**

$$R_D(f) = \frac{1}{D} \sum_D \text{loss}(y_i, f(x_i))$$

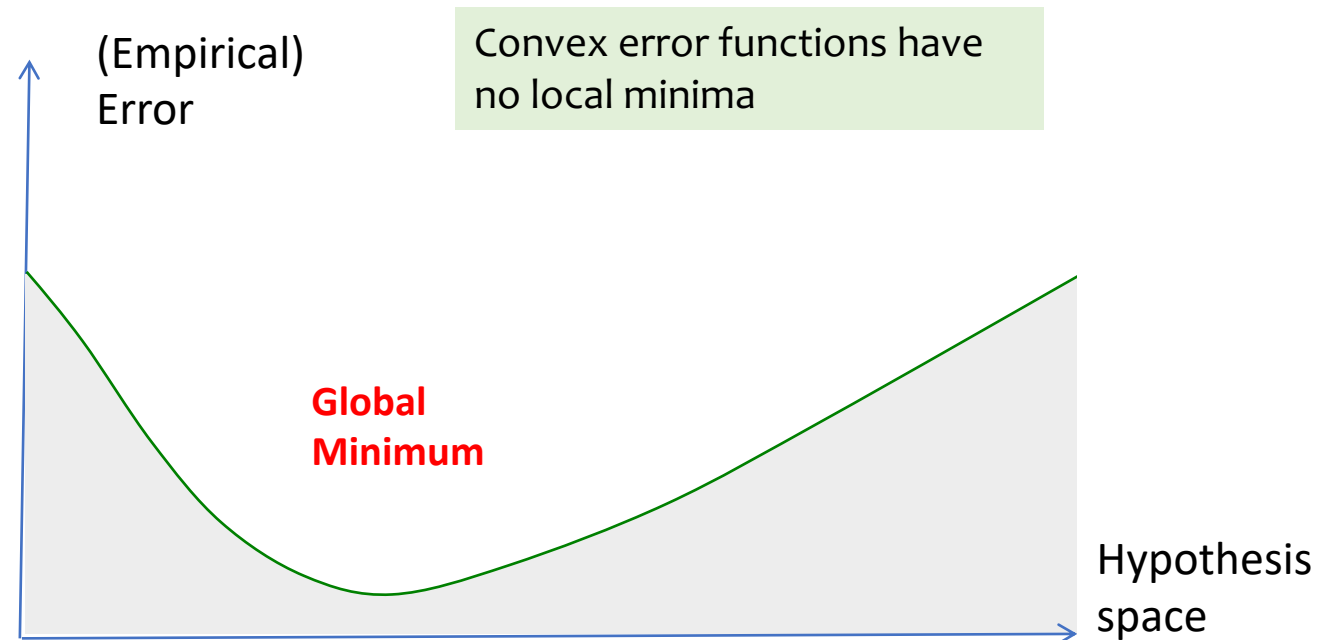
# Error Surface

- **Linear classifiers:** hypothesis space parameterized by  $w$
- *Error/Loss/Risk* are all functions of  $w$



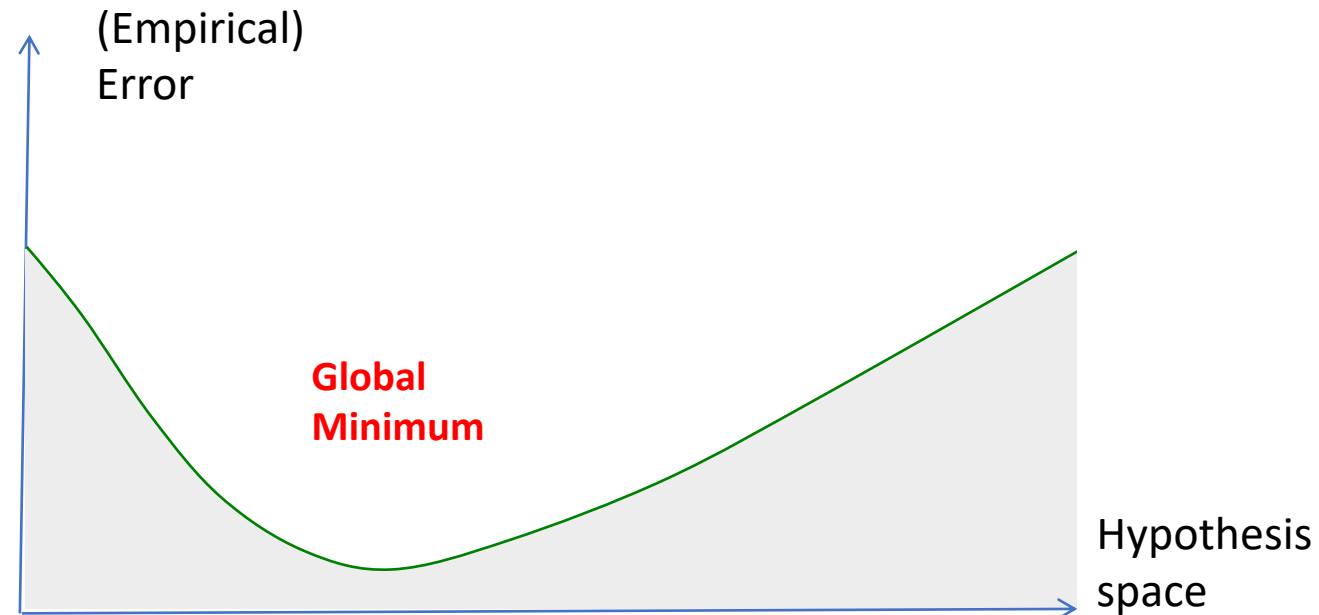
# Convex Error Surfaces

- **Convex functions** have a single minimum point
  - Local minimum = global minimum
  - *Easier to optimize*



# Loss minimization

- **Let's consider the square loss.**
  - Convex loss function, error surface has a global minimum (~any local minimum is also global).



# Gradient Descent for Squared Loss

$$Err(w) = \frac{1}{2} \sum_{d \in D} (y_d - f(x_d))^2$$

Initialize  $\mathbf{w}^0$  randomly  
for  $i = 0 \dots T$ :

$\Delta \mathbf{w} = (0, \dots, 0)$

    for every training item  $d = 1 \dots D$ :

$f(\mathbf{x}_d) = \mathbf{w}^i \cdot \mathbf{x}_d$

        for every component of  $\mathbf{w}$       $j = 0 \dots N$ :

$\Delta w_j += \alpha(y_d - f(\mathbf{x}_d)) \cdot x_{dj}$

$\mathbf{w}^{i+1} = \mathbf{w}^i + \Delta \mathbf{w}$

return  $\mathbf{w}^{i+1}$  when it has converged

# Loss minimization

- **Let's consider the square loss.**
  - Convex loss function, error surface has a global minimum (~any local minimum is also global).





# Expected vs. Empirical Risk

- The risk (aka generalization error) of a classifier is its **expected loss** (the loss averaged over all possible datasets)

$$R(f) = \int L(y, f(x))P(x, y)dx, y$$

$$\epsilon \triangleq \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(y, f(x))] = \sum_{(x,y)} \mathcal{D}(x, y) \ell(y, f(x))$$

$$R_D(f) = \frac{1}{D} \sum_D \text{loss}(y_i, f(x_i))$$

The empirical risk of a classifier on a dataset D is its average loss on the items in d

# Loss minimization

- **Let's consider the square loss.**

- Convex loss function, error surface has a global minimum (~any local minimum is also global).



*Do we really want to get to that global minimum point?*

- *we care about minimizing the expected loss, while this error surface describes the empirical loss!*

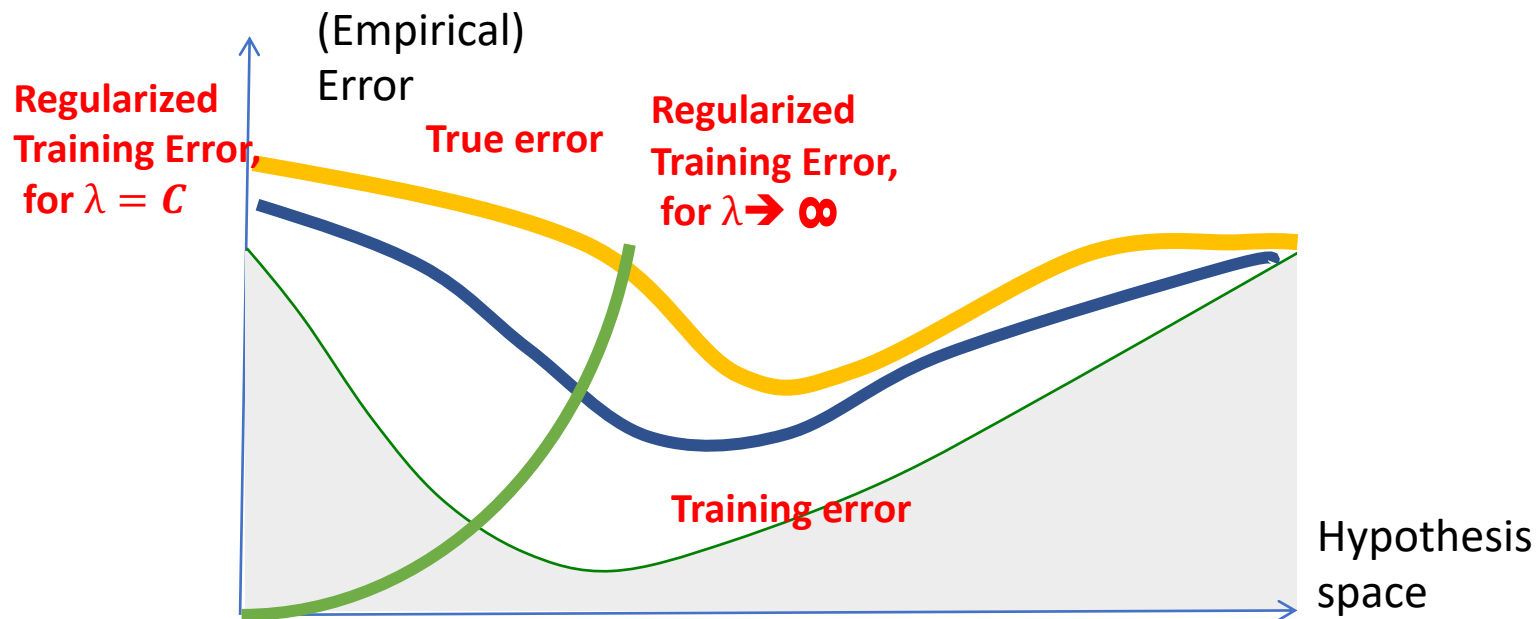
# Regularization

- *A form of inductive bias – we prefer simpler functions!*
- A very popular choice of regularization term is to minimize the norm of the weight vector
  - For convenience:  $\frac{1}{2}$  squared norm

$$\min_{\mathbf{w}} \sum_n \text{loss}(y_n, \mathbf{w}_n) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

# Loss minimization

- **Let's consider the square loss.**
  - Convex loss function, error surface has a global minimum (~any local minimum is also global).

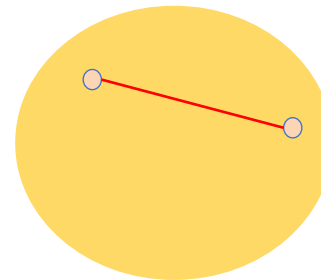
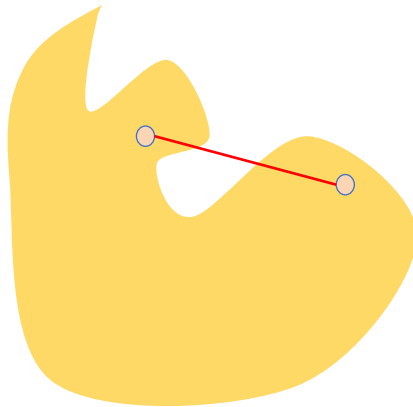


# Convex Sets

- A **set**  $S$  is **convex** if

$$\forall x, y \in S, \alpha x + (1 - \alpha)y \in S \quad (0 \leq \alpha \leq 1)$$

*(the line segment joining  $x$  and  $y$  is contained in  $S$ )*



# Convex Function

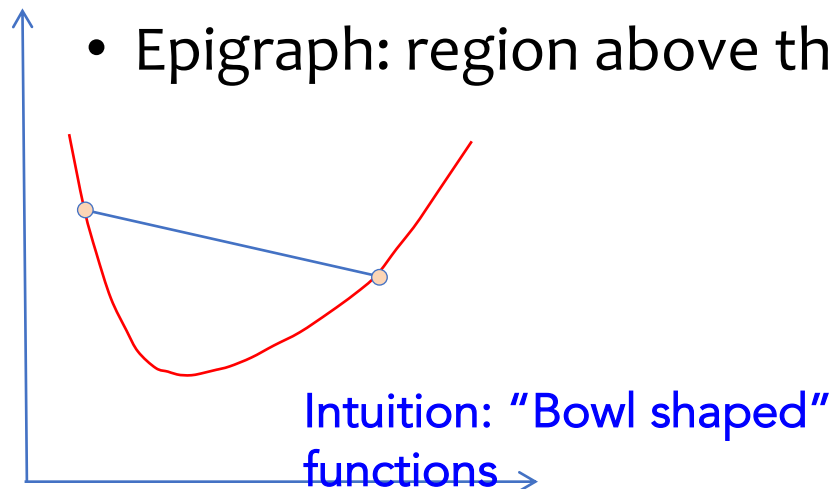
- A function  $f$  defined on a convex set is convex if

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \quad (0 < \alpha < 1)$$

$f$  is convex if the chord joining any two points is always above the graph.

- A function  $f$  is convex if its epigraph is a convex set

- Epigraph: region above the graph of the function  $f$



If  $f$  is convex  $\rightarrow -f$  is concave

# Checking Convexity

- According to definition: chord lies above function

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \quad (0 < \alpha < 1)$$

- If  $f$  is differentiable:  $f$  lies above all tangent lines

$$f(y) \geq f(x) + f'(x)(y - x) \quad f(y) \geq f(x) + \nabla f(x)(y - x)$$

- For twice differentiable functions: 2<sup>nd</sup> derivative is non-negative

$$f''(x) \geq 0 \quad \nabla^2 f(x) \succeq 0$$

- If the above functions are strict inequalities,  $f$  is strictly convex

# Example: Checking Convexity

$$f(x) = x \log x$$

$$f'(x) = \log x + 1$$

$$f''(x) = \frac{1}{x}$$

- $f''(x) > 0$  for all  $x > 0 \Rightarrow f(x)$  is (strictly) convex!
  - *Note that the domain of  $\log(x)$  is  $x > 0$*

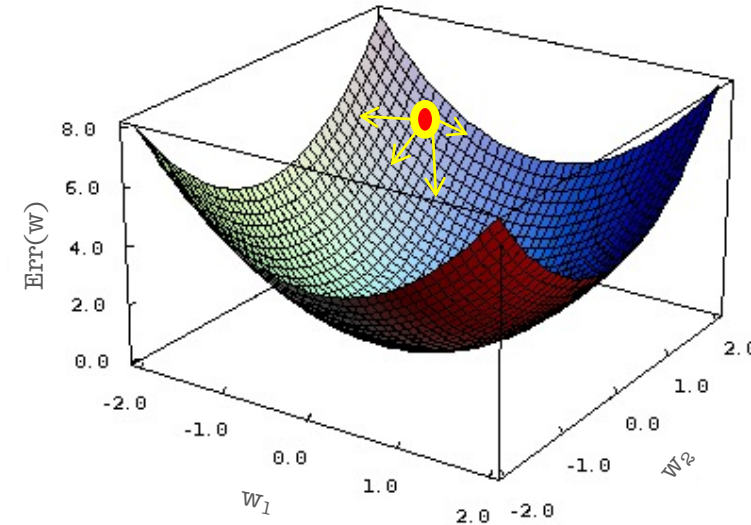


# Error Surface for Squared Loss

We add the  $\frac{1}{2}$  for convenience

$$Err(w) = \frac{1}{2} \sum_{d \in D} (\hat{y}_d - y_d)^2$$

$y = w^T x$



Since  $\hat{y}$  is a constant (for a given dataset), the Error function is a *quadratic function* of  $W$  (paraboloid)

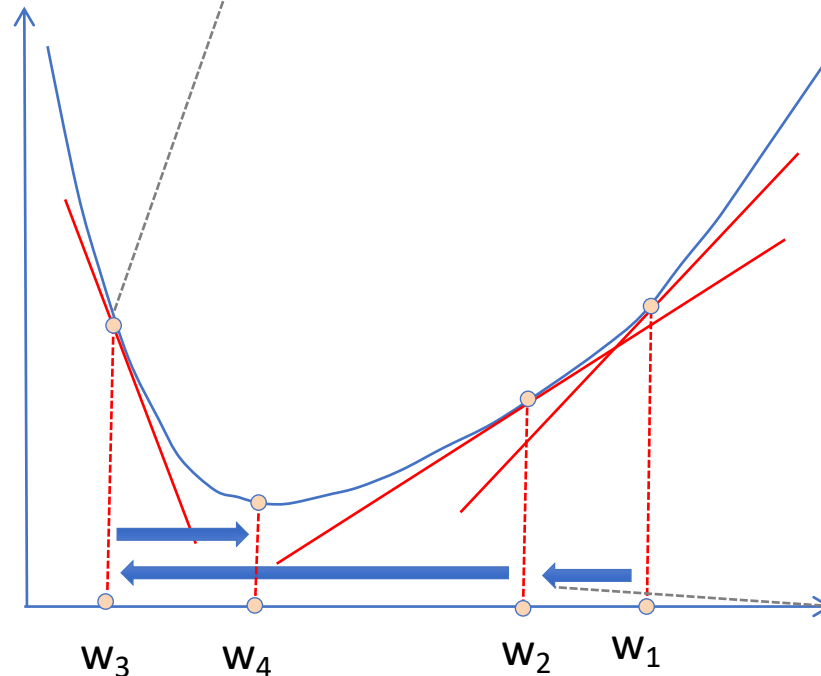
➔ Squared Loss function is convex!

How can we find the global minimum?

# Gradient Descent Intuition

(3) We also need to determine the step size (aka learning rate).

*What happens if we overshoot?*



(1) The derivative of the function at  $w_1$  is the slope of the tangent line  
→ Positive slope (increasing)

*Which direction should we move to decrease the value of  $Err(w)$ ?*

(2) The gradient determines the direction of steepest increase of  $Err(w)$  (go in the opposite direction)

*What is the gradient of  $Error(w)$  at this point?*

# The Gradient of Error(w)

The gradient is a generalization of the derivative

$$\nabla Err(w) = \left( \frac{\partial Err(w)}{\partial w_0}, \frac{\partial Err(w)}{\partial w_1}, \dots, \frac{\partial Err(w)}{\partial w_n} \right)$$

The gradient is a vector of partial derivatives.

It Indicates the *direction of steepest increase* in  $Err(w)$ , for each one of  $w$ 's coordinates

# Computing $\nabla \text{Err}(\mathbf{w}^i)$ for Squared Loss

$$\frac{\partial \text{Err}(\mathbf{w})}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (y_d - f(\mathbf{x}_d))^2$$

$$\text{Err}(w) = \frac{1}{2} \sum_{d \in D} (y_d - f(x_d))^2$$

$$= \frac{1}{2} \frac{\partial}{\partial w_i} \sum_{d \in D} (y_d - f(\mathbf{x}_d))^2$$

$$= \frac{1}{2} \sum_{d \in D} 2(y_d - f(\mathbf{x}_d)) \frac{\partial}{\partial w_i} (y_d - \mathbf{w} \cdot \mathbf{x}_d)$$

$$= - \sum_{d \in D} (y_d - f(\mathbf{x}_d)) x_{di}$$

# Batch Update Rule for Each $w_i$

- **Implementing gradient descent:** *as you go through the training data, accumulate the change in each  $w_i$  of  $W$*

$$\Delta w_i = \alpha \sum_{d=1}^D (y_d - \mathbf{w}^i \cdot \mathbf{x}_d) x_{di}$$

# Gradient Descent for Squared Loss

```
Initialize  $\mathbf{w}^0$  randomly
for  $i = 0 \dots T$ :
     $\Delta \mathbf{w} = (0, \dots, 0)$ 
    for every training item  $d = 1 \dots D$ :
         $f(\mathbf{x}_d) = \mathbf{w}^i \cdot \mathbf{x}_d$ 
        for every component of  $\mathbf{w}$   $j = 0 \dots N$ :
             $\Delta w_j += \alpha(y_d - f(\mathbf{x}_d)) \cdot x_{dj}$ 
     $\mathbf{w}^{i+1} = \mathbf{w}^i + \Delta \mathbf{w}$ 
return  $\mathbf{w}^{i+1}$  when it has converged
```

# Batch vs. Online Learning

- The Gradient Descent algorithm makes updates after going over the entire data set
  - Data set can be huge
  - Streaming mode (we cannot assume we saw all the data)
  - Online learning allows “adapting” to changes in the target function
- Stochastic Gradient Descent
  - Similar to GD, updates after each example
  - Can we make the same convergence assumptions as in GD?
- Variations: update after a subset of examples

# Stochastic Gradient Descent

Initialize  $\mathbf{w}^0$  randomly

for  $m = 0 \dots M$ :

$$f(\mathbf{x}_m) = \mathbf{w}^i \cdot \mathbf{x}_m$$

$$\Delta w_j = \alpha(y_d - f(\mathbf{x}_m)) \cdot x_{mj}$$

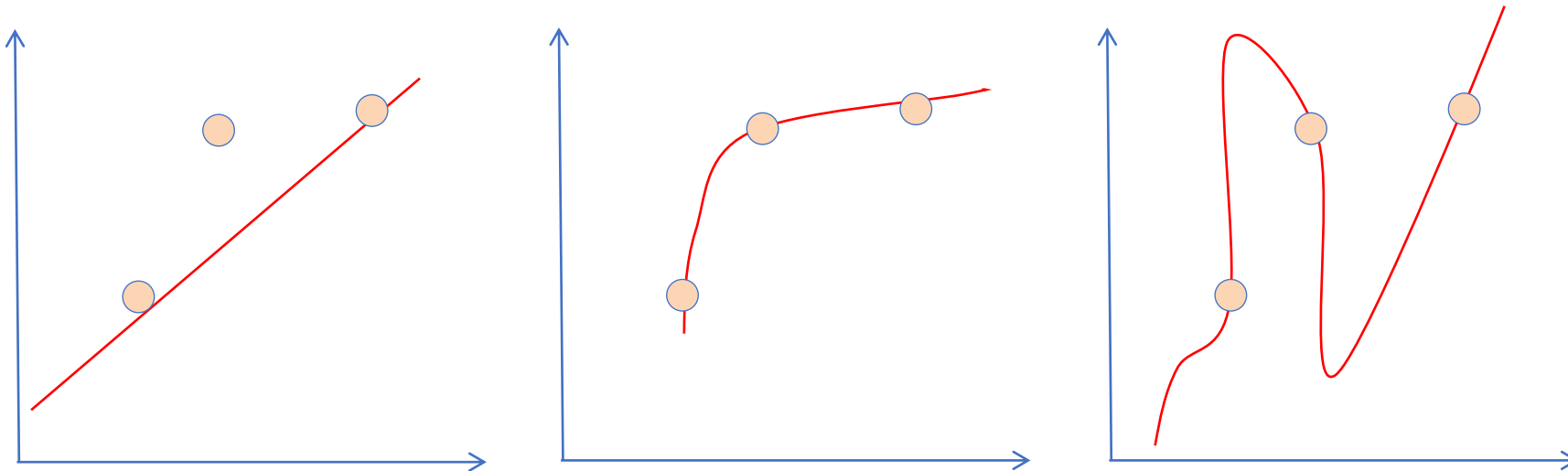
$$\mathbf{w}^{i+1} = \mathbf{w}^i + \Delta \mathbf{w}$$

return  $\mathbf{w}^{i+1}$  when it has converged



# Polynomial Regression

- GD: general optimization algorithm
  - Works for classification (different loss function)
  - Incorporate polynomial features to fit a complex model
  - **Danger – overfitting!**

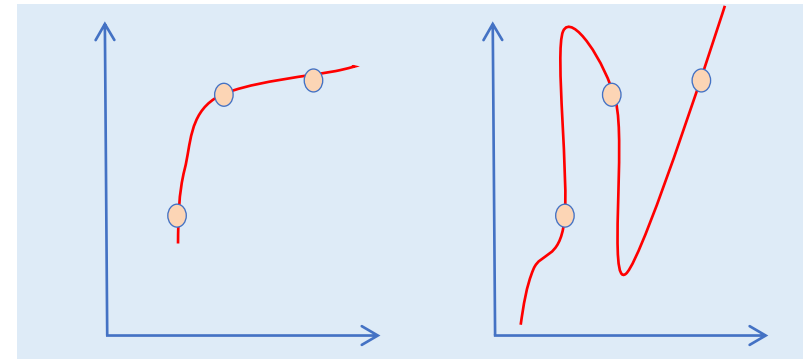
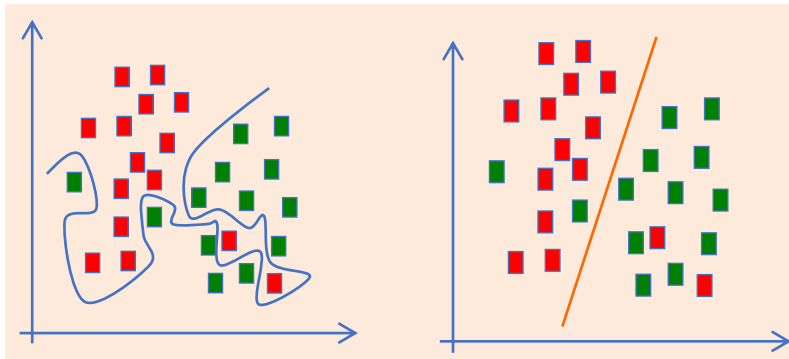


# Regularization

- Both for regression and classification, for a given error we prefer a *simpler model*
  - *Keep  $W$  small:  $\epsilon$  changes in the input cause  $\epsilon * w$  in the output*
- Some times we are even willing to trade a higher error rate for a simpler model (*why?*)
- *Add a regularization term:*
  - This is a form of **inductive bias**

$$\min_w = \sum_n \text{loss}(y_n, w_n) + \lambda R(w)$$

*How different values affect learning?*



# Regularization

- A very popular choice of regularization term is to minimize the norm of the weight vector

- For example  $||\mathbf{w}|| = \sqrt{\sum_d w_d^2}$

- For convenience:  $\frac{1}{2}$  squared norm

$$\min_{\mathbf{w}} = \sum_n \text{loss}(y_n, \mathbf{w}_n) + \frac{\lambda}{2} ||\mathbf{w}||^2$$

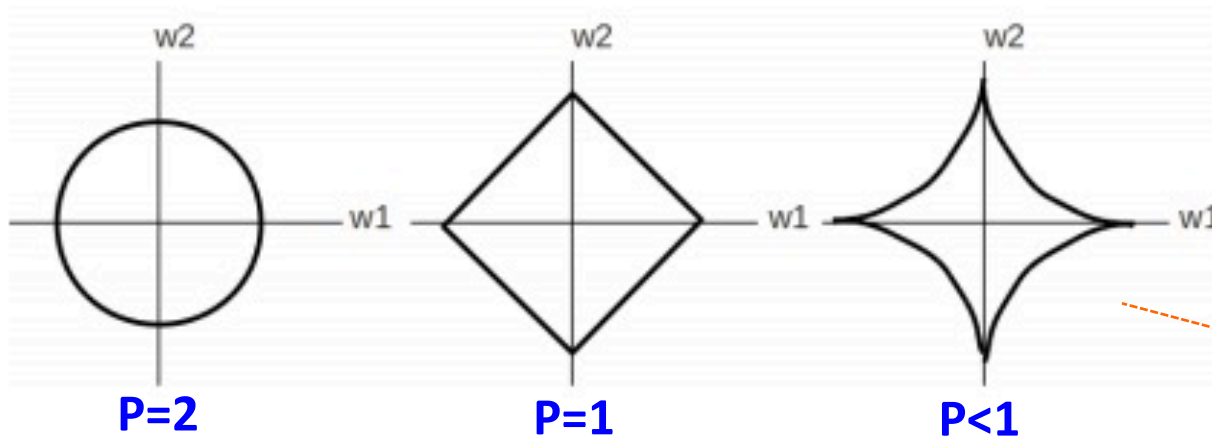
- *What is the update gradient of the loss function?*
  - *At each iteration we subtract the weights by  $\lambda * \mathbf{w}$*
- In general, we can pick other norms (p-norm)
  - Referenced as *L-p norm*
  - Different p will have a different effect!
  - *What will the norm for p=1 and p=0 minimize?*

$$||\mathbf{w}||_p = \left( \sum_d |w_d|^p \right)^{\frac{1}{p}}$$

# Different Regularization functions

- To understand the difference between different norms consider their iso-surface  $||w||_p = C$  (\*c is a constant value)

$$||w||_p = \left( \sum_d |w_d|^p \right)^{\frac{1}{p}}$$



**p<1**: peaked: large values of one weight comes at the "expense" of another

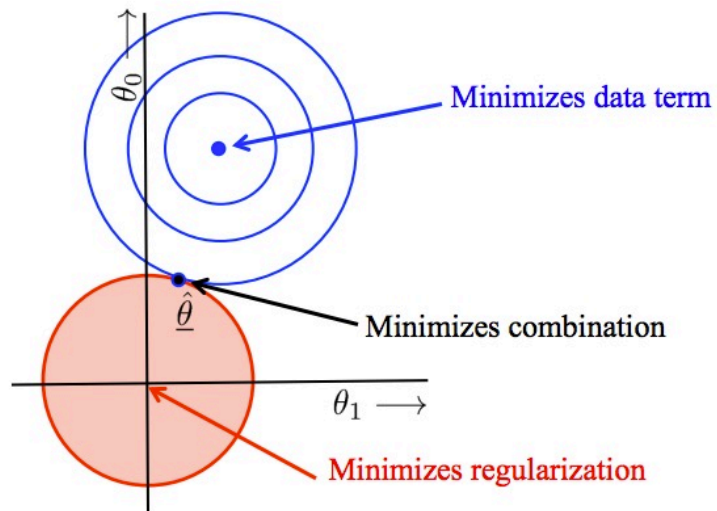
- As the value of p gets closer to 0, the norm is closer to a **counting norm**: count non-zero parameters (ignoring values) → directly minimizes the number of active features
- In general: smaller values of p encourage sparsity

How would the iso surface of p=0 look?

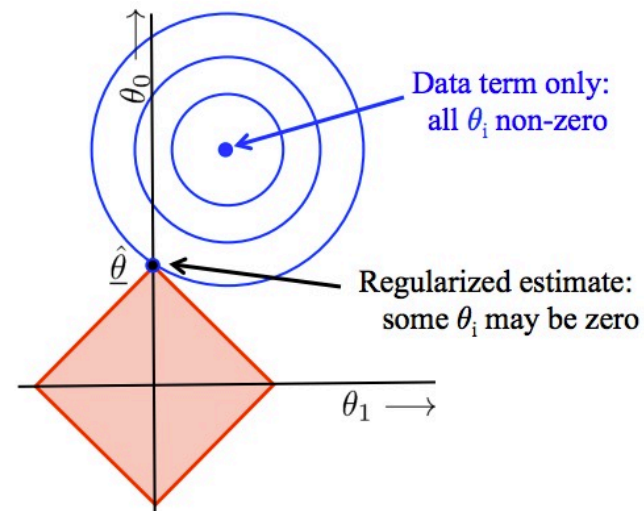
# L1 and L2 norms

- Common choices: L1 and L2 are both convex ( $L_{p<1}$  is not convex)
- *Regularized objective*: balance between minimizing the error and the regularization cost

*L2 optimum will be sparse,  
ONLY if the data loss term is  
minimized at the axis*



*L1 norm contour are **sharp**, will intersect with the contour of data loss term even when the data loss term min point is not at the axis  
→ L1 encourages sparsity*



# Classification

- **So far:**

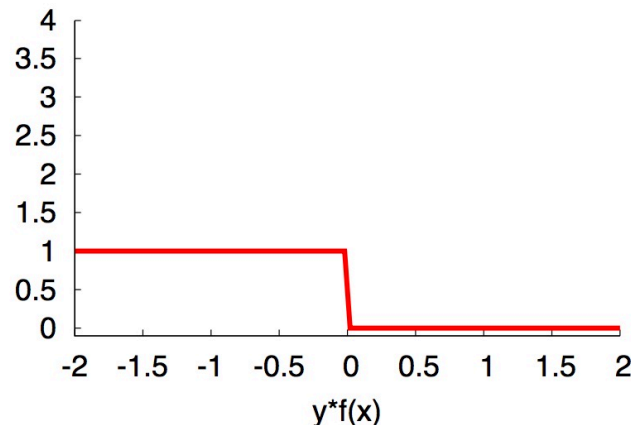
- General **optimization framework for learning**
- Minimize regularized loss function

$$\min_{\mathbf{w}} = \sum_n \text{loss}(y_n, \mathbf{w}_n) + \frac{\lambda}{2} ||\mathbf{w}||^2$$

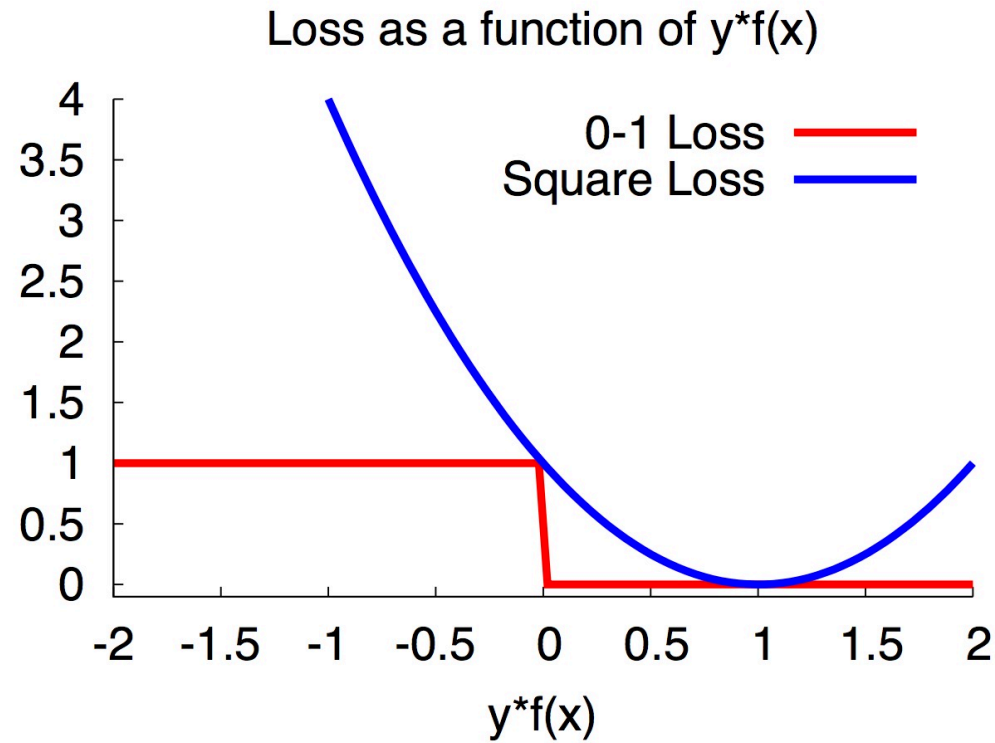
- Gradient descent is an all purpose tool
  - Computed the gradient of the **square loss function**
- **Moving to classification should be very easy**
  - Simply replace the loss function
- ...

# Classification

- Can we minimize the empirical error directly?
  - Use the 0-1 loss function
- **Problem:** *Cannot be optimized directly*
  - *Non convex and non differentiable*
- **Solution:** define a smooth loss function, an upper bound to the 0-1 loss function



# Square Loss is an Upper bound to 0/1 Loss

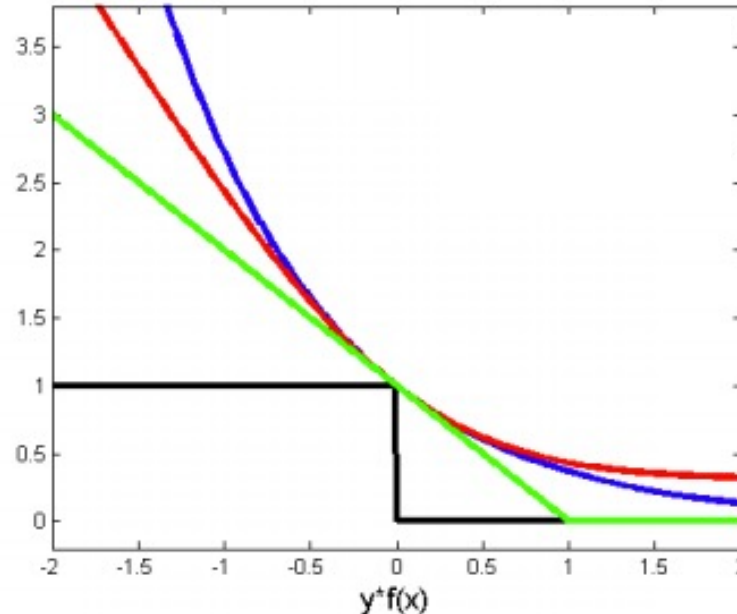


Is the square hinge loss a good candidate to be a surrogate loss function for 0-1 loss?



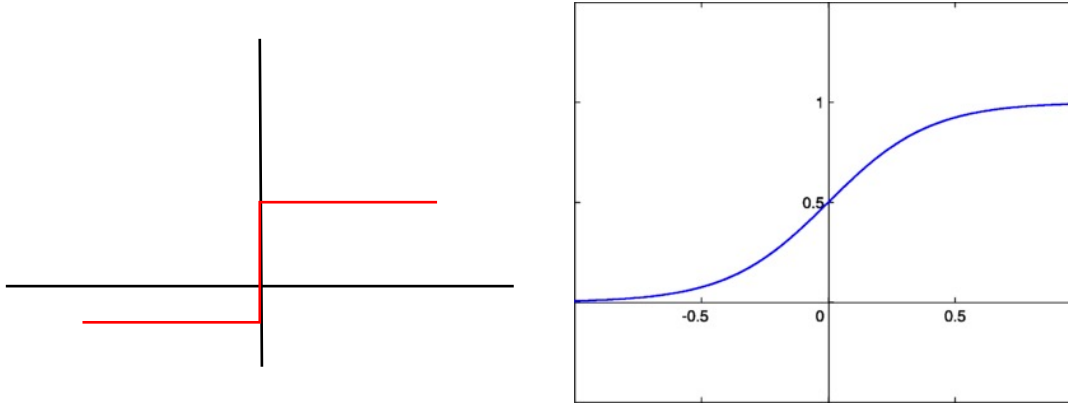
# Surrogate Loss functions

- Surrogate loss function: smooth approximation to the 0-1 loss
  - Upper bound to 0-1 loss



# Logistic Function

- Smooth version of the threshold function



$$h_w(x) = g(w^T x)$$

$$z = w^T x$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

↑ Sigmoid (logistic) function

- Known as a sigmoid/logistic function
  - Smooth transition between 0-1
- Can be interpreted as the conditional probability
- Decision Boundary
  - $y=1: h(x) > 0.5 \rightarrow w^T x \geq 0$
  - $y=0: h(x) < 0.5 \rightarrow w^T x < 0$

# Logistic Regression

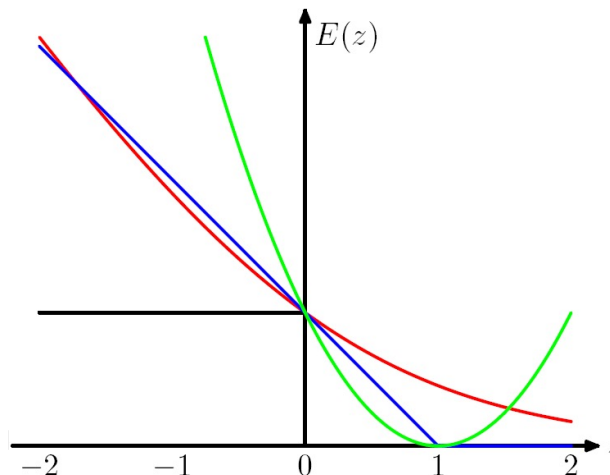
- **Learning:** optimize the likelihood of the data
  - **Likelihood:** *probability of our data under current parameters*
  - For easier optimization, we look into the log likelihood (*negative*)
- **Cost Function**
  - If  $y = 1$ :  $-\log P(y=1 | x, w) = -\log (g(w, x))$ 
    - *If the model gives very high probability to  $y=1$ , what is the cost?*
  - If  $y = 0$ :  $-\log P(y=0 | x, w) = -\log (1 - P(y=1 | x, w)) = -\log (1 - g(w, x))$
- Or more succinctly (with  $l_2$  regularization)

$$Err(w) = - \sum_i y^i \log(e^{g(w, x_i)}) + (1 - y^i) \log(1 - e^{g(w, x_i)}) + \frac{1}{2} \lambda ||w||^2$$

- This function is convex and differentiable
  - Compute the gradient, minimize using gradient descent

# Question

- 0-1 loss has a steep change when  $y(w^T x) > 0$  and  $y(w^T x) < 0$ , while the logistic loss function has a smooth transition.
- This means that even correct classification will incur some cost.
- *How does this fact bias the learning process?*

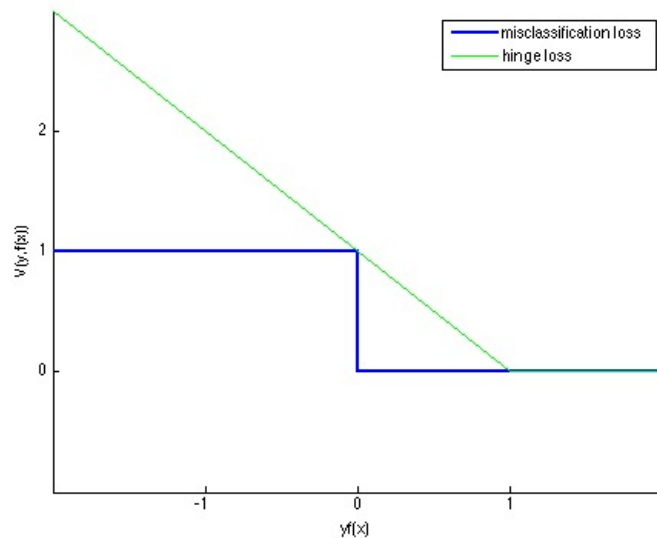


# Hinge Loss

- Another popular choice for loss function is the hinge loss

$$L(y, f(x)) = \max(0, 1 - y f(x))$$

- We will discuss in the context of support vector machines (SVM)



It's easy to observe that:

- (1) The hinge loss is an upper bound to the 0-1 loss
- (2) The hinge loss is a good approximation for the 0-1 loss
- (3) BUT ...

It is not differentiable at  $y(w^T x) = 1$

**Solution:** Sub-gradient descent

# Summary

- Introduced an optimization framework for learning:
  - Minimization problem
  - Objective: data loss term and regularization cost term
  - Separate learning objective from learning algorithm
    - Many algorithms for minimizing a function
- Can be used for regression and classification
  - Different loss function
  - GD and SGD algorithms
- Classification: use surrogate loss function
  - Smooth approximation to 0-1 loss