

Introduction to Machine Learning

Lecture 9

Instructor:
Dr. Tom Arodz



Recap: the story so far...

- Ideal way of making predictions:
 - Perfect feature
 - Accurate, efficient simulation
 - Fully-known probability distribution $p(y|x)$ from which data comes from
 - or $p(x|y)$, we can use Bayes theorem to get $p(y|x)$ from it
- Traditional machine learning:
 - heuristics / ad-hoc approaches (e.g. perceptron in HW1)
 - maximum likelihood (or Bayesian learning)
 - pick a distribution shape for $p(x|y)$, based on domain knowledge
 - fit best distribution $p(x|y)$ to each class in training data
 - get $p(y|x)$ from these distributions using Bayes theorem
- “Modern” machine learning:
 - produce numbers that we can treat as $p(y|x)$ from some function $f(x)$ – not any specific distribution
 - Tweak the parameters θ of function $f(x; \theta)$ to make predictions better

Towards ML Approach #3

- “Modern” machine learning: probability distrib. governing the problem not known, let's model it using a class of function
 - Predicting: use trained parameters θ_{opt}
 - $p(\text{class}|\mathbf{x}, S) = \text{some_f}(\theta_{\text{opt}})(\mathbf{x}).\text{to_class_probs}()$
 - Training:
 - Search for single value of parameters θ that minimizes:
$$\theta_{\text{opt}} = \arg \min_{\theta} \text{error_metric}(S, \text{some_f}(\theta)(\mathbf{x} \text{ in } S).\text{to_class_probs}())$$
- Compare with maximum likelihood (MLE) :
 - Predicting: use trained parameters θ_{opt}
 - $p(\text{class}|\mathbf{x}, S) = \text{mult_normal}(\theta_{\text{class}}).\text{pdf}(\mathbf{x}) * P[\text{class}]$
 - Training:
 - $\theta_{\text{class}} = \arg \max_{\theta} \text{mult_normal}(\theta).\text{pdf}(S_{\text{class}})$

Recap: ML Approach #3

■ “Modern” machine learning:

```
minimizeθ: error_metric(S,  
    some_f(θ)(x in S).to_class_probs())
```

■ Filling the gaps:

- `some_f(θ)(x)` or `some_f(θ;x)`
 - Some function that takes two groups of inputs
 - Raw input features x
 - Trainable parameters θ
 - the parameters help us pick a specific function from a family of functions
- `to_class_probs()`
 - `some_f(θ;x)` may not return probabilities, how to convert to: ≥ 0 , $\text{sum}=1$
- `error_metric`
 - Classification error? Do we need something more complex?
- `minimizeθ`
 - how? For Gaussians in MLE, we had analytical solution, that's unlikely for `some_f`



9.1. Iterative minimization



Building ML Approach #3

- Modern machine learning:

```
minimize $\theta$ : error_metric(  
    S, some_f( $\theta$ ) (x in S).to_class_probs())
```

- Minimization:

- minimize θ
 - how?

Building ML Approach #3

- Modern machine learning:

```
minimizeθ: error_metric(  
    S, some_f(θ)(x in S).to_class_probs())
```

- Let's simplify the notation:

- minimize_θ:
$$\sum_{(x,y) \in S} \text{error_metric}(\quad$$
$$y,$$
$$\text{some_f}(\theta, x).to_class_probs())$$

- minimize_θ: $\sum_{(x,y) \in S} L(y, \text{model}(\theta, x))$

Building ML Approach #3

- Modern machine learning:
 - minimize_{θ} :
$$\sum_{(x,y) \in S} \text{error_metric}(y, \text{some_f}(\theta, x).to_class_probs())$$
- Let's assume that `error_metric`, `some_f`, and `to_class_probs` are smooth functions: they have derivatives (as opposed to e.g. $|x|$ - abs. value of x does not have derivative at $x=0$)
- That means: `L` and `model` are smooth functions, too (have derivatives)
 - minimize_{θ} : $\sum_{(x,y) \in S} L(y, \text{model}(\theta, x))$
 - Note that Σ (sum) is also a smooth function
- Problems to solve: find minimum of a differentiable function
 - minimize_{θ} : `function(θ)`

Building ML Approach #3

■ Training:

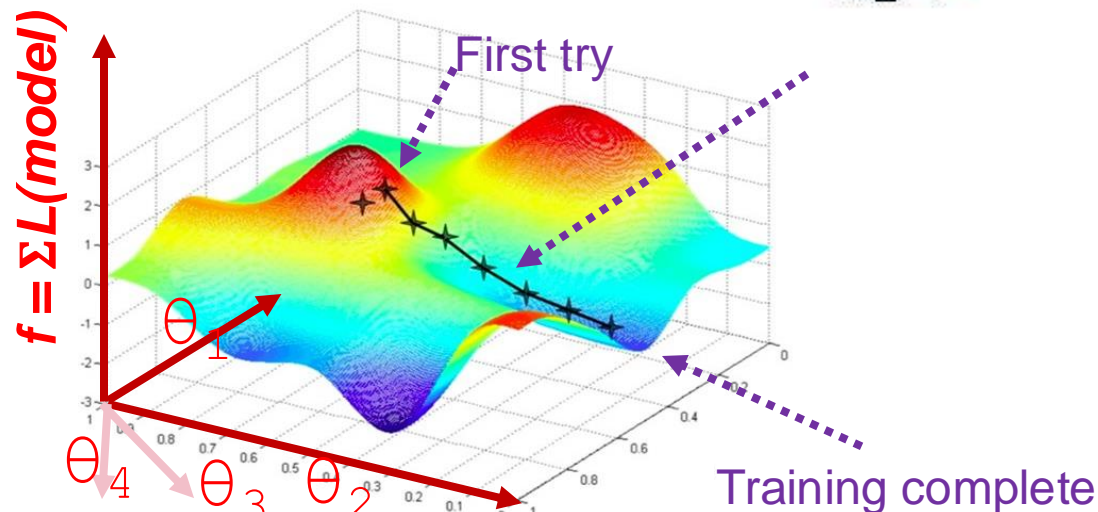
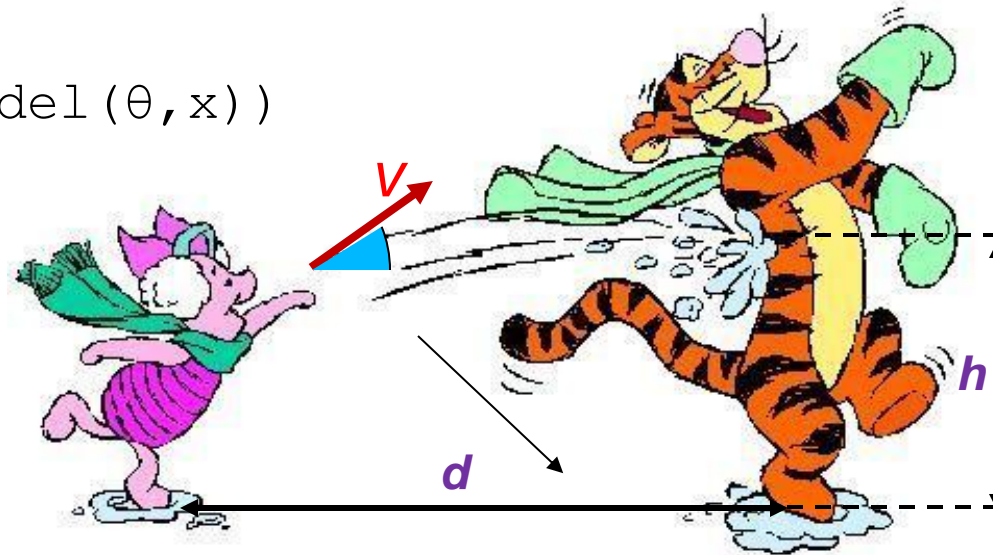
- $\text{minimize}_{\theta}: \text{function}(\theta)$
 - More specifically:
- $\text{minimize}_{\theta}: \sum_{(x,y) \in S} L(y, \text{model}(\theta, x))$

■ Like Perceptron / HW1

but

we assumed
differentiability of \mathbf{f}

it will help us figure out
in which direction
we should tweak θ





Optimization

- Training:
 - $\text{minimize}_{\theta}: \text{function}(\theta)$
- We will change notation to just minimize:
 - $\text{minimize}_x: f(x)$
- For a moment:
 - x will not mean “feature values” of our training set
 - instead, it will mean the parameters we want to find to minimize the value of $f()$
 - the training data will be “folded in” (hidden inside) the shape of the function $f()$
 - that means: $f(x)$ is very complex, its shape depends on training set, and difficult to view/analyze “holistically”



Optimization

- Task

- $\text{minimize}_x: f(x)$

- Possible ways to find minimum of a function:

- Analytical: solve for "derivative of $f = 0$ "

- E.g. $\min_x x^2 - 4x$ is at $d(x^2 - 4x)/dx = 2x - 4 = 0$, i.e. at $x = 2$

- We used this approach in Maximum Likelihood Estimation for Gaussians

- Cons:

- in ML, f is highly complex, finding the equation to solve and solving it will be computationally hard
 - in ML, f has many points with $df/dx = 0$: maxima, saddle points, many minima

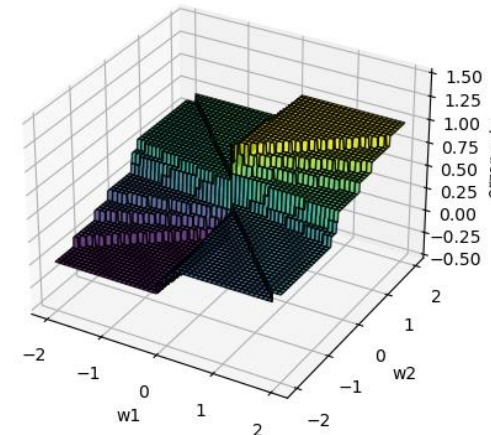
Optimization

■ Task

- $\text{minimize}_x: f(x)$

■ Possible ways to find minimum of a function:

- Analytical won't work, ultimately because:
 - $f(x)$ is very complex, its shape depends on training set, and difficult to view/analyze "holistically"
- What is "cheap" is calculating value of the function $f(x)$ for specific x . So, one alternative to "analytical" is:
- Probing: like constructing the plots in HW1 with meshgrid
 - try $f(x)$ at a large number of possible x
 - pick best, i.e., lowest value of $f(x)$



Optimization

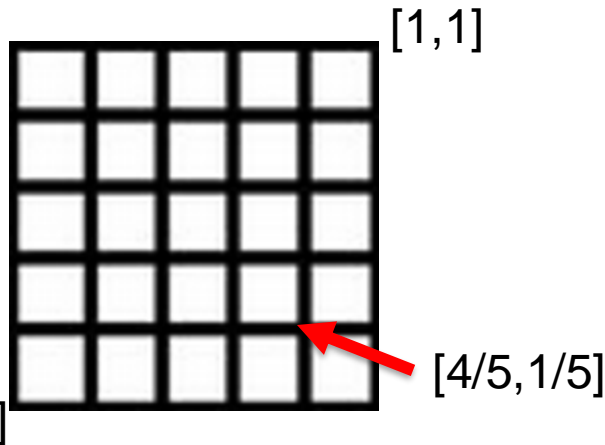
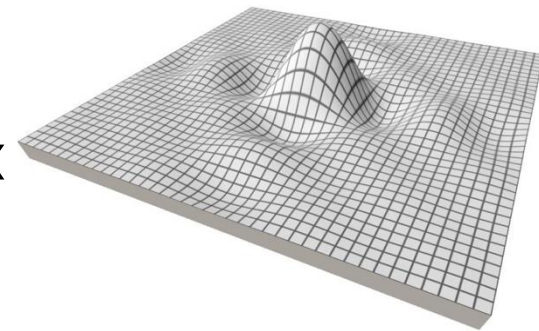
■ Task

- $\text{minimize}_x: f(x)$

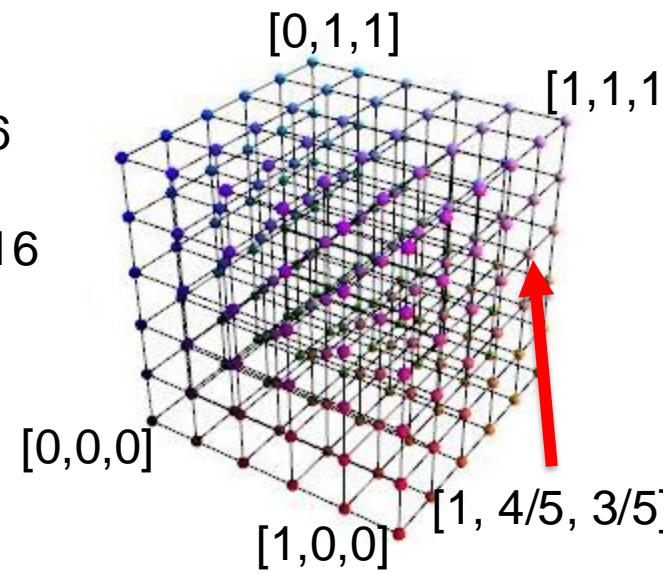
■ Possible ways to find minimum of a function:

■ Probing: like plots in HW1

- try $f(x)$ at a large number of possible x
- pick best, i.e., lowest value of $f(x)$



- $p=5$,
- If $n=2$, we need $6^2=36$ points in the grid
- If $n=3$, we need $6^3=216$ points in the grid
- each cell has side length of $1/5$



- Won't work, to many points to try!

Optimization

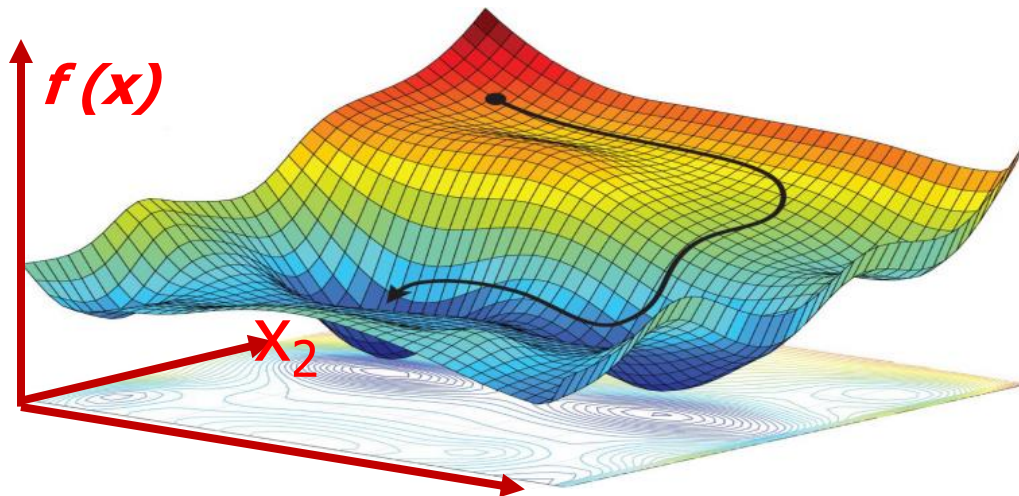
- Task

- $\text{minimize}_x: f(x)$

- Possible ways to find minimum of a function:

- Iterative minimization

- Like training in HW1



- Relies on “local” rules to find how to update in position (x here, w in HW1)

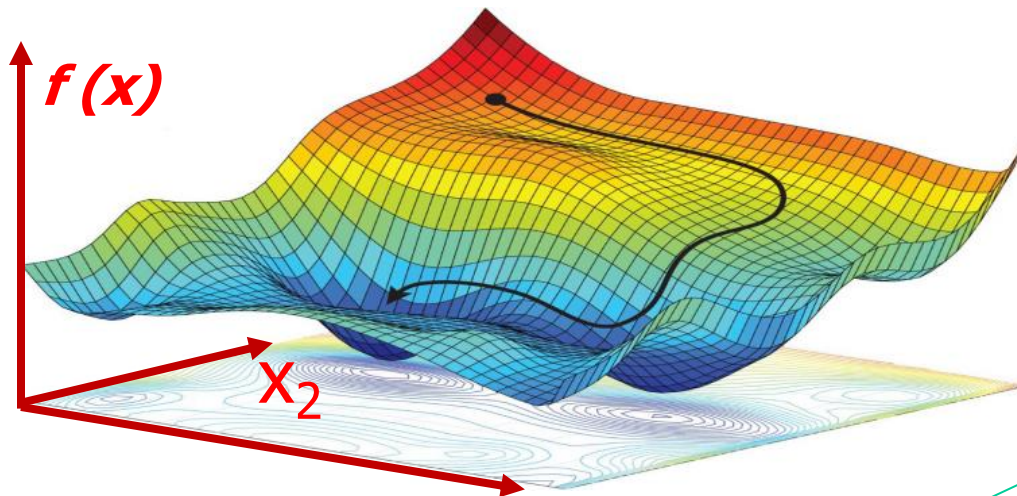
Optimization

- Task

- $\text{minimize}_x: f(x)$

- Possible ways to find minimum of a function:

- Iterative minimization



- Relies on “local” rules to find how to update in position
 - Heuristic rules like in HW1
 - Rules based on value of $f(x)$ for neighboring x
 - Rules based on derivatives (gradient) of f at x

ML overwhelmingly uses this approach: “gradient descent”

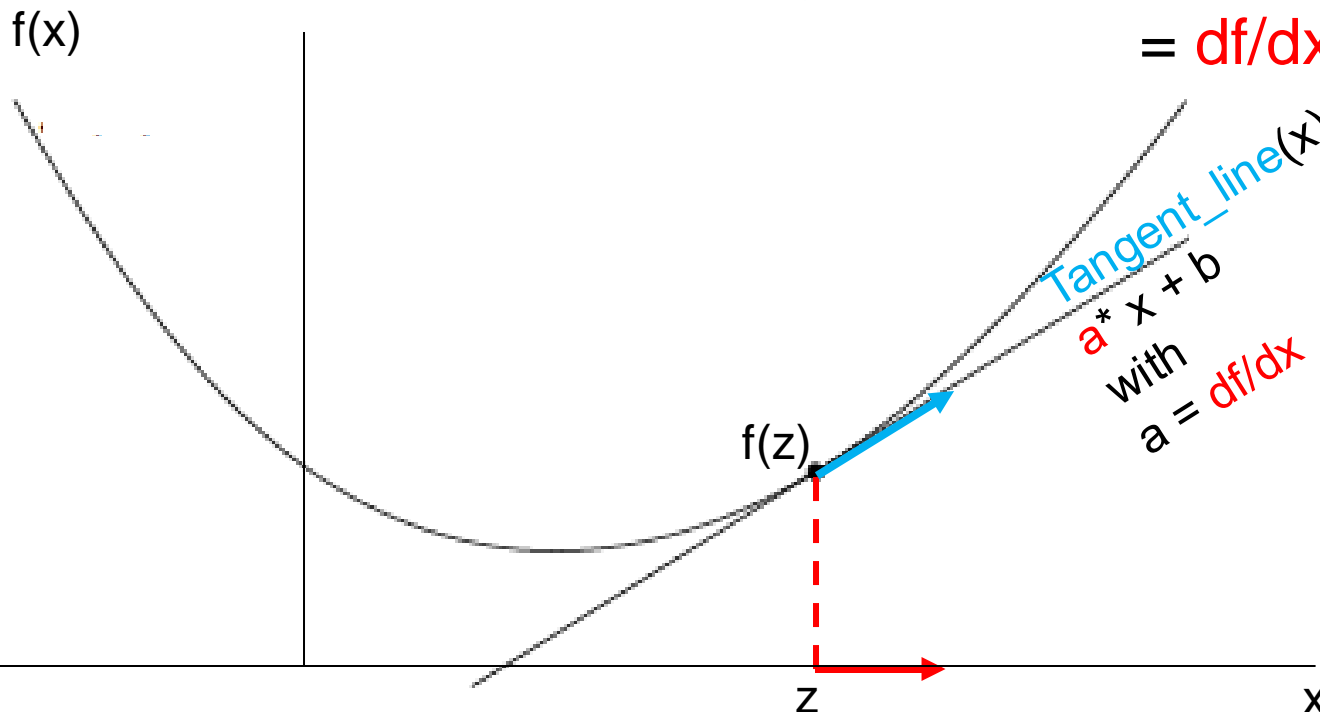
Gradient descent

For 1D, $\nabla f = df/dx = f'$

First derivative or gradient: denoted ∇f or f'

gradient of a n-dimensional function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a generalization of derivative of a single-dimensional function $f: \mathbb{R} \rightarrow \mathbb{R}$

$$\begin{aligned} \text{Tangent_line}(x) &= f(z) + df/dx * (x-z) \\ &= df/dx * x + [f(z) - df/dx * z] \end{aligned}$$



Gradient descent

gradient of a n-dimensional function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is:

an n-dimensional **vector** of partial derivatives,

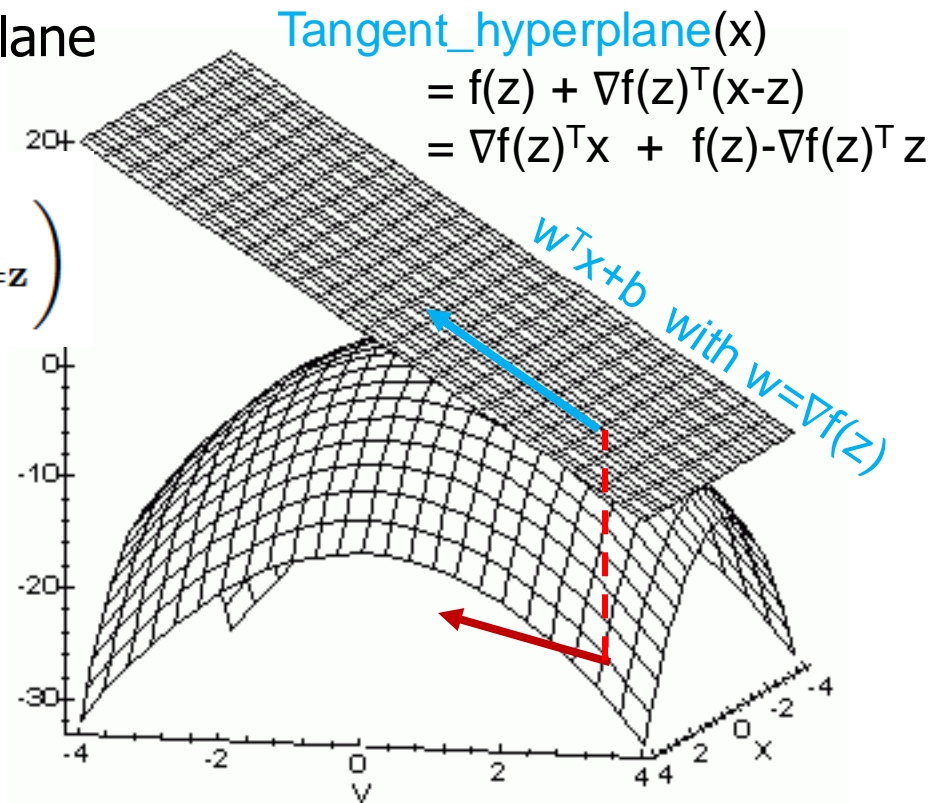
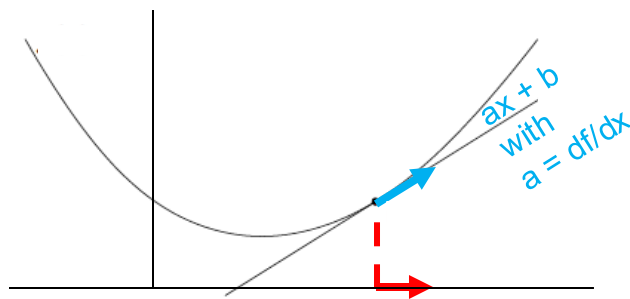
pointing (for a given z) in direction of steepest slope of f at that z

Pointing always in the direction in which the function grows (locally)

positive value in a vector = f grows (vector points) towards $+\infty$; negative value = f grows towards $-\infty$, vector points towards $-\infty$

Defines a line or plane or hyperplane tangent to f

$$\nabla f(\mathbf{z}) = \left(\frac{\partial}{\partial x_1} f(\mathbf{x})|_{\mathbf{x}=\mathbf{z}}, \dots, \frac{\partial}{\partial x_n} f(\mathbf{x})|_{\mathbf{x}=\mathbf{z}} \right)$$



Gradient descent

Gradient descent (with learning rate c):

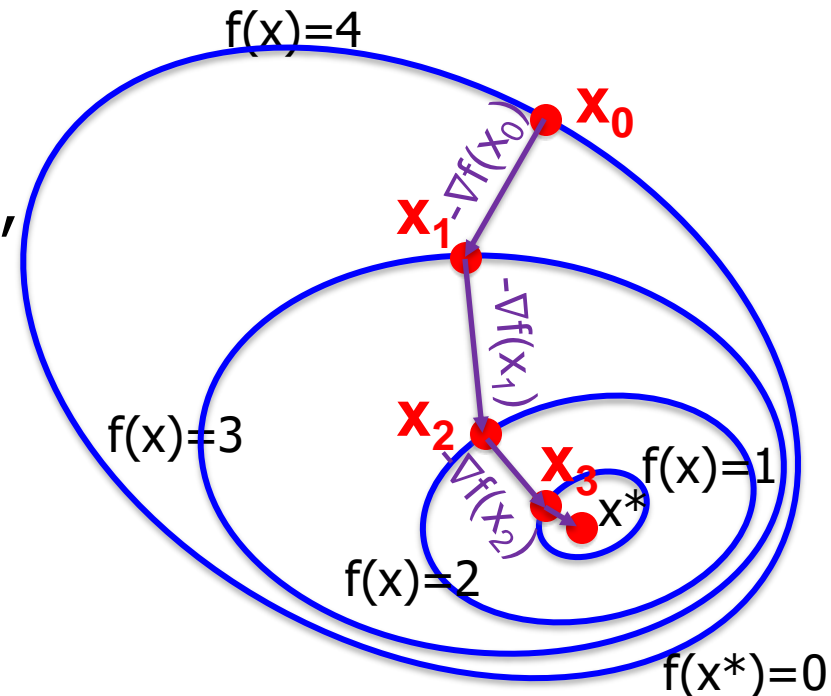
We start from x_0

We calculate $x_1 = x_0 - c \nabla f(x_0)$

We calculate $x_2 = x_1 - c \nabla f(x_1)$

$x_{n+1} = x_n - c \nabla f(x_n)$

If we choose L large enough
grad.desc. goes down in each step,
converging towards
local minimum of function f

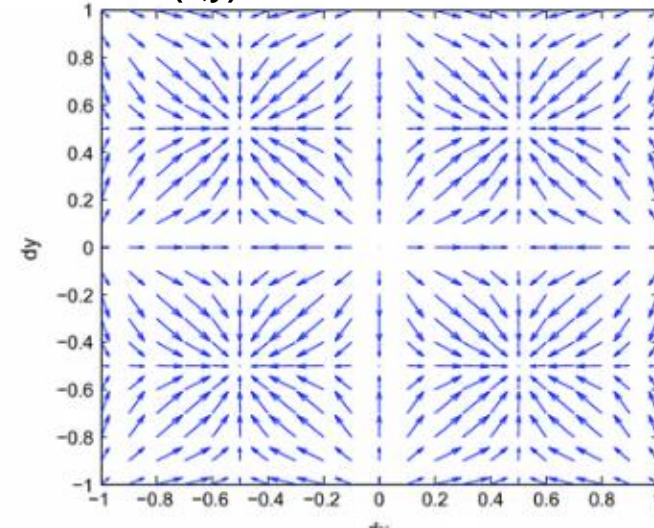
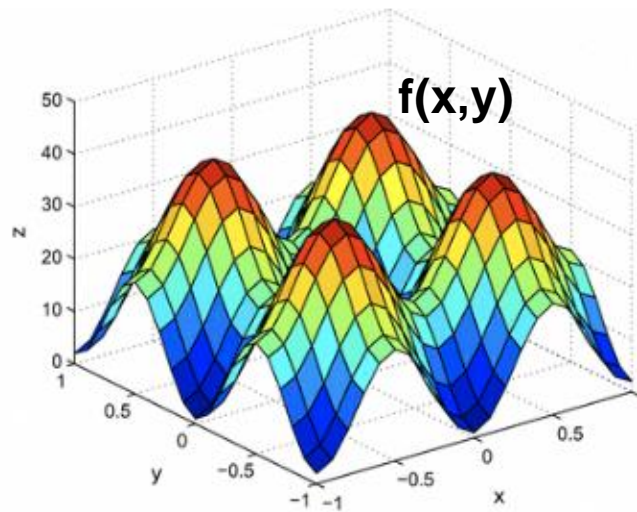


Recap: Gradient descent

$$\nabla_{(x,y)} f = [\partial f / \partial x, \partial f / \partial y]$$

Gradient of $f(\mathbf{x})$ at \mathbf{z} :

vector of partial derivatives of \mathbf{f} w.r.t. \mathbf{x} at the current point \mathbf{z} pointing towards the direction of steepest increase of function locally at point \mathbf{z}



Gradient descent (with learning rate c):

We start from \mathbf{x}_0

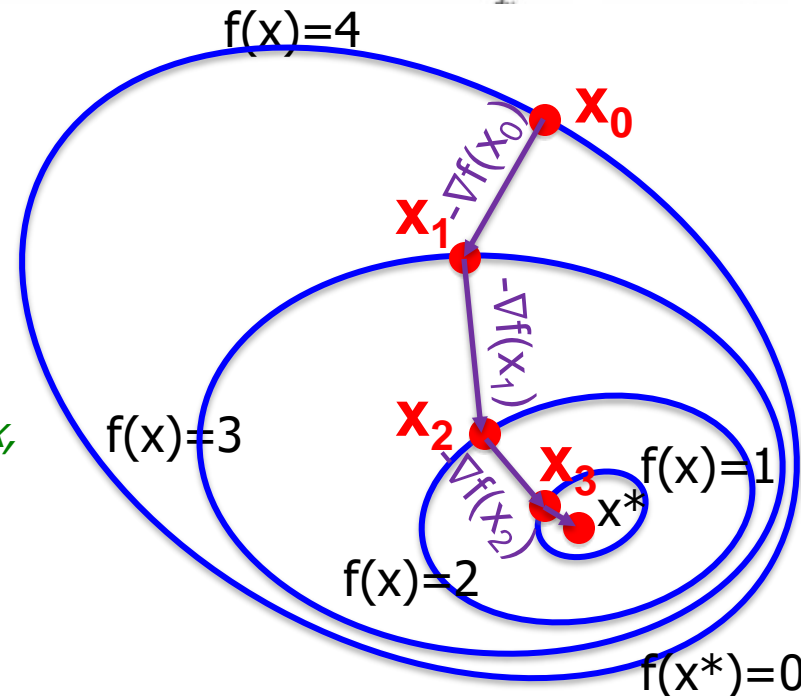
We calculate $\mathbf{x}_1 = \mathbf{x}_0 - c \nabla f(\mathbf{x}_0)$

We calculate $\mathbf{x}_2 = \mathbf{x}_1 - c \nabla f(\mathbf{x}_1)$

$\mathbf{x}_{n+1} = \mathbf{x}_n - c \nabla f(\mathbf{x}_n)$

Note: we don't need full formula of ∇f as function of x , just the code to get its value for specific x_0, x_1, \dots

In ML, we apply grad. desc. to find minimum of a function of parameters θ , (e.g. $\theta = (w, b)$), instead of function of x



Recap

■ Modern machine learning:

```
minimizeθ: error_metric(  
    S, some_f(θ) (x in S).to_class_probs())
```

■ Filling the gaps:

- `some_f(θ) (x)` **or** `some_f(θ; x)`
 - Some function (a “model”) that takes two groups of inputs
 - Raw input features x
 - **Trainable parameters θ**
- `to_class_probs()`
 - `some_f(θ; x)` may not return probabilities, how to convert to: ≥ 0 , $\text{sum}=1$
- `error_metric`
 - Classification error? Do we need something more complex?
- `minimizeθ`
 - **How?**
Gradient descent, over parameters θ of the model `some_f(θ; x)`
 - **$\theta_{n+1} = \theta_n - c \nabla_{\theta} \text{error_metric}(S, \text{model}(\theta_n, S))$**



9.2. Correct/incorrect as the *error metric*

Building ML Approach #3

$\text{minimize}_{\theta}: \text{error_metric}(S, \text{model}(\theta, S))$

- Assume “model” return raw predictions from some_f, or predictions converted to class probabilities

■ Filling the gaps:

■ `error_metric`

- Classification error (error rate or error count) is the most natural metric of classifier performance
 - $\text{error_rate}(\text{training_set } S) = \text{error_count}(\text{training_set } S) / |S|$
 - $\text{error_count}(\text{training_set } S) = \sum_{(x, y) \in S} \text{is_incorrect}(y, \text{predicted_y}(x))$
 - `is_incorrect` returns 1 (incorrect) or 0 (correct)
- With gradient descent as the minimization method, the update of the parameters θ would be (c =learning rate):
 - $\text{new } \theta = \text{old } \theta - c \nabla_{\theta} \text{error_count}(S, \text{model}(\theta, S))$
or, equivalently
 - $\text{new } \theta = \text{old } \theta - c \sum_{(x, y) \in S} \nabla_{\theta} \text{is_incorrect}(y, \text{model}(\theta, x))$

Choosing the error metric

error_metric

- Classification error (error rate or error count) is the most natural metric of classifier performance

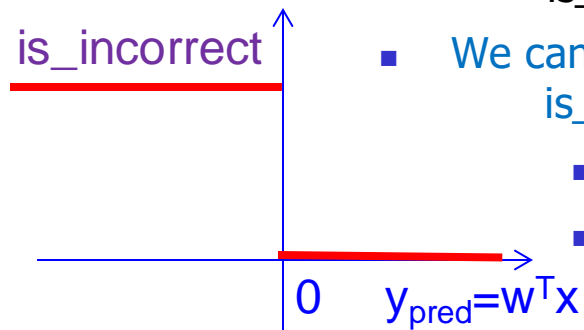
- $\text{error_rate}(\text{training_set } S) = \text{error_count}(\text{training_set } S) / |S|$
- $\text{error_count}(\text{training_set } S) = \sum_{(x, y) \in S} \text{is_incorrect}(y, \text{predicted_y}(x))$
 - is_incorrect returns 1 (incorrect) or 0 (correct)

- We can also use:

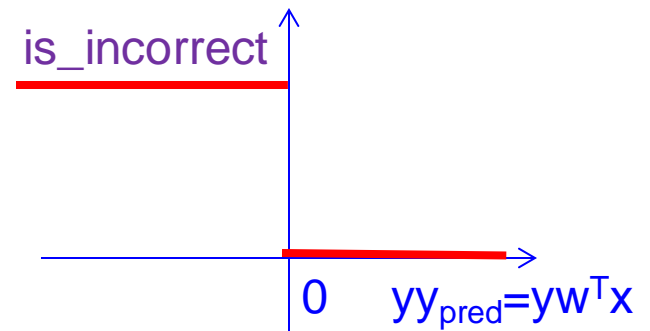
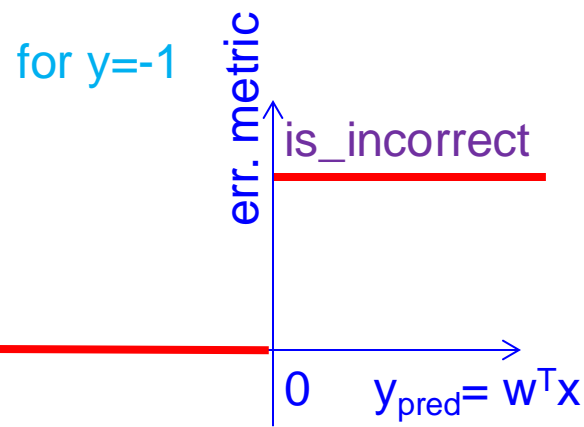
$$\text{is_incorrect}(y, \text{predicted_y}(x)) = \text{is_negative}(y * \text{predicted_y}(x))$$

- $\text{is_negative}(u)$ returns 1 for $u \leq 0$, and returns 0 otherwise
- can be coded as $\text{sign}(1 + \text{sign}(-1 * y * \text{predicted_y}(x)))$

for $y=+1$



for $y=-1$



when using the two-part term
 $y_k f(w; x_k) = y_{\text{true}} * y_{\text{pred}}$
 there's just one common plot

Choosing the error metric

```
minimizeθ: error_metric(  
    S, some_f(θ) (x in S).to_class_probs())
```

- Filling the gaps:

- minimize_θ
 - how? Gradient descent!
- error_metric
 - Classification error is the most natural metric of classifier performance



How to use it
with gradient descent?

∇ error_count = ?

∇ is_incorrect = ?

∇ is_negative($yw^T x$) = ?

Choosing the error metric

```
minimizeθ: error_metric(  
    S, some_f(θ) (x in S).to_class_probs())
```

■ Filling the gaps:

- minimize_θ
 - how? **Gradient descent!**
- error_metric
 - Classification error is the most natural metric of classifier performance



$\nabla \text{error_count} = 0$

$\nabla \text{is_incorrect} = 0$

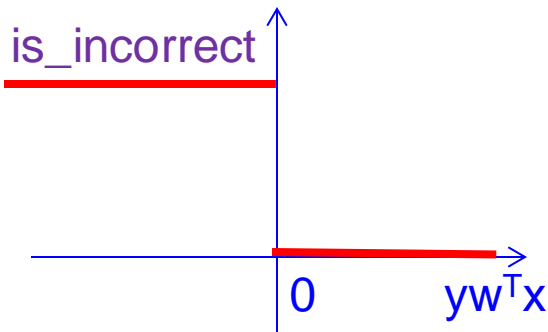
$\nabla_{\text{is_negative}}(\mathbf{y}^T \mathbf{x}) = 0$

Flat everywhere ($\nabla=0$)
except
when

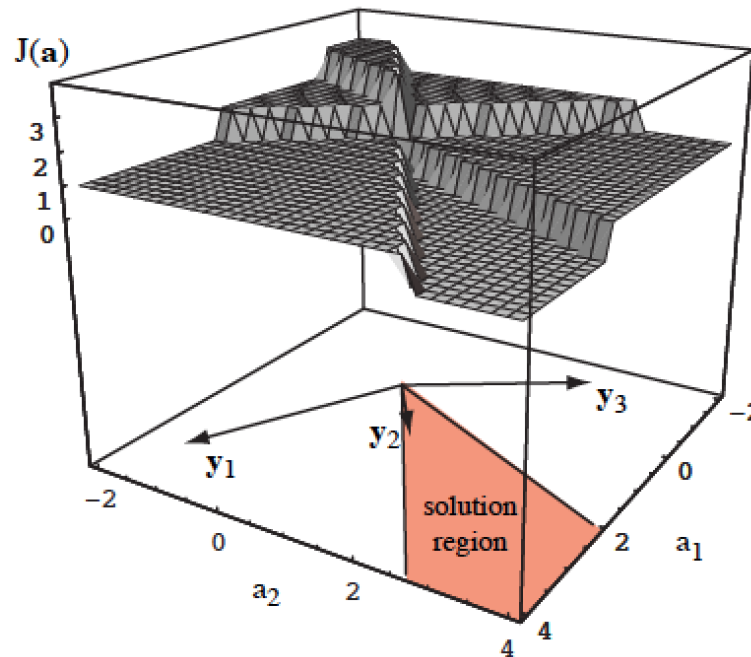
is_incorrect changes
from 1 to 0,
error_count drops by 1
(sudden drops:
derivative not defined)

Choosing the error metric

- minimize_{θ}
 - how? Gradient descent!
- `error_metric`
 - Classification error is the most natural metric of classifier performance



0 if $yh(x) > 0$,
1 otherwise



$\nabla \text{error_count}$

$\nabla \text{is_incorrect}$

Non-differentiable

Not continuous

Flat!

- We don't know in which direction to move to get a (even slightly) better solution!



Choosing the error metric

```
minimizeθ: error_metric(  
    S, some_f(θ)(x in S).to_class_probs())
```

- Filling the gaps:

- `error_metric`

- Classification error is the most natural metric of classifier performance
 - But will not work with gradient-based (or other local) optimization approaches
 - We will need a different error metric
 - One that has no flat regions!



9.3. Mean squared error



Towards ML Approach #3

`minimizeθ: error_metric(S, model(θ, S))`

■ Filling the gaps:

- `minimizeθ`

- how? **Gradient descent!**

- `error_metric`

- **Mean square error** is a well-known error metric:

- $\text{MSE}(\text{training_set } S) = \sum_{(x,y) \in S} (y - \text{predicted_y}(x))^2 / |S|$

- We assume that true y is +1 or -1,

- With gradient descent as the minimization method,
the update of the parameters θ would be (c =learning rate):

- $\text{new } \theta = \text{old } \theta - c \nabla_{\theta} \text{MSE}(S, \text{model}(\theta, S))$

or, equivalently

- $\text{new } \theta = \text{old } \theta - c/|S| \sum_{(x,y) \in S} \nabla_{\theta} (y - \text{model}(\theta, x))^2$

Is MSE a good choice for classification?

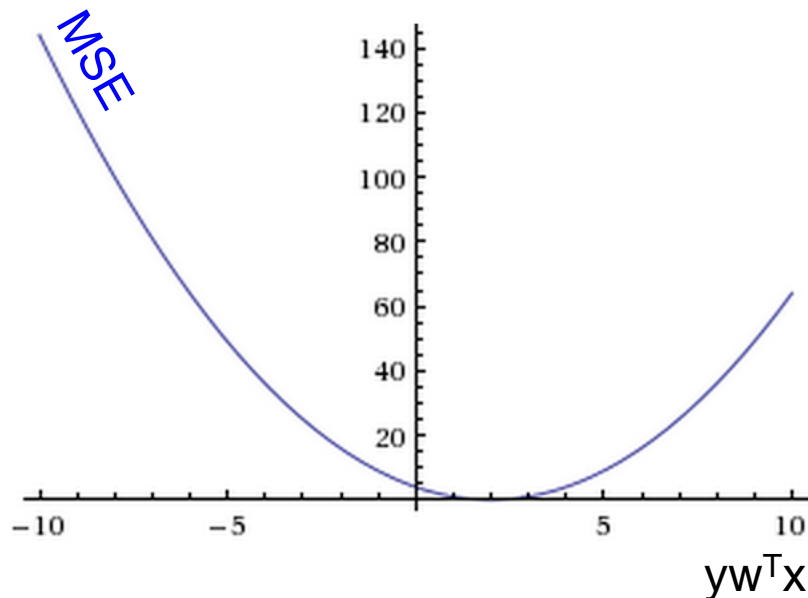
■ MSE for classification:

$$\text{loss}(w, x, y) = (y - w^T x)^2 \quad \text{multiply by } y^2=1$$

~~$= (1 - yw^T x)^2$~~

■ MSE is:

- Differentiable (gradient exists everywhere)
- Not flat (we will get non-zero gradient)
- Nonnegative



Is MSE a good choice for classification?

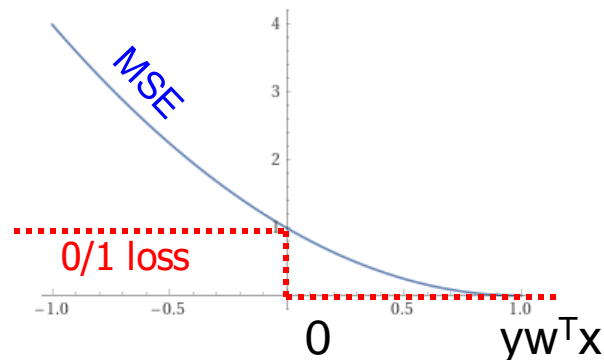
■ MSE for classification:

$$\text{loss}(w, x, y) = (y - w^T x)^2 \quad \text{multiply by } y^2 = 1$$

■ MSE:

- Has some relationship to 0/1 loss (is-incorrect)
 - If MSE = 0 then 0/1_loss = 0
 - Upper-bounds 0/1 loss:

$$\text{MSE}(w, x, y) \geq \text{0/1_loss}(w, x, y)$$



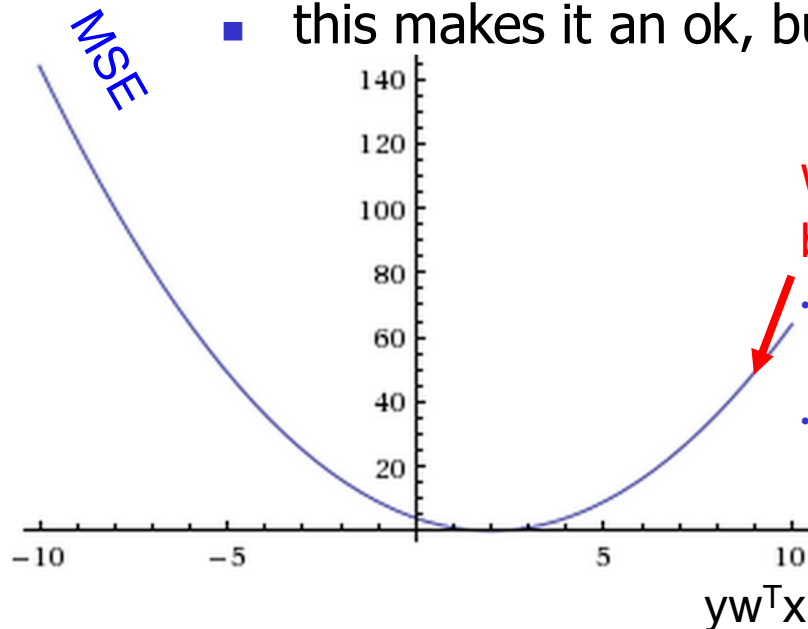
Problem with MSE for classification

- MSE is:

- Not flat
- Upper-bounds 0/1 loss

- Non-monotonic:

- decreases, then increases, as $yw^T x$ increases
 - E.g. for class +1, prediction +4 (correct) has same MSE as prediction -2 (incorrect)
- this makes it an ok, but not the best choice for classification



We get large loss (penalty) not only for incorrect, but for correct predictions too!

- If data is Gaussian, very few points should be that far from decision boundary.
- But what if data is not Gaussian? BAD!
 - non-monotonic loss is suspicious for classification (but ok for regression)



Updated schedule

- Test 1 will be on Wednesday, 10/9
 - Test will be done via Canvas, online (from home)
- Study problems will be posted in Canvas on Friday, 9/27
- We will discuss the study problems during our Monday, 10/7 class
- HW2 is due next Tuesday, 10/1, at 5pm