

# MATH 565: Lecture 1 (01/13/2026)

1-1

Today: \* logistics, syllabus, ...  
\* problems in ML, optimization for them

This is Optimization for Machine Learning (Math 565)  
I'm Bala Krishnamoorthy (call me Bala).

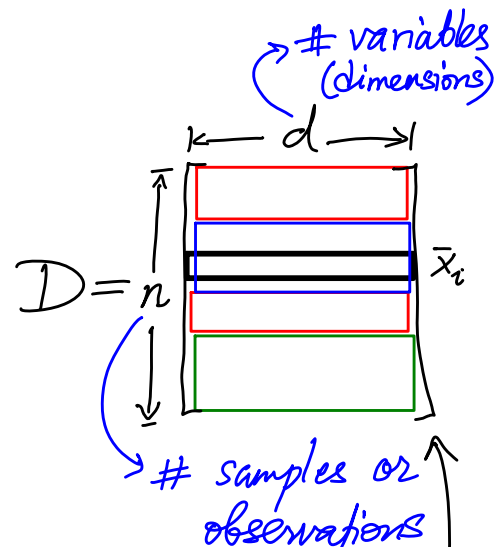
We will **not** use Canvas. All materials for the class will be posted on the course web page: <https://bala-krishnamoorthy.github.io/Math565.html>

Download the Book PDF (accessible via WSU Libraries)!

- \* Check the syllabus on the course web page.
- \* homework assignments will be posted at least a week before its due.

## Machine Learning Problems

1. **Clustering** We are given a data matrix that is  $n \times d$  for  $n$  observations (or samples) each having values for  $d$  variables (dimension= $d$ ). The  $i$ th observation is denoted  $\bar{x}_i$  ( $d$ -vector), which forms the  $i$ th row of  $D$  (when written as  $\bar{x}_i^T$ ).



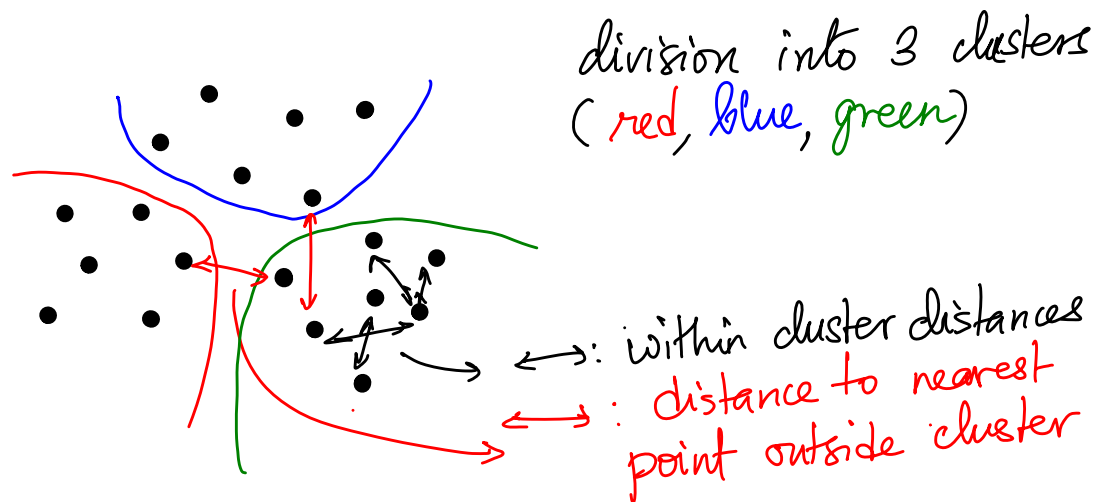
The goal of clustering is to divide the  $n$  observations into  $K$  disjoint subsets or clusters. The # clusters  $K$  may be user-defined, e.g., in  $k$ -means clustering, or may be determined as part of the clustering method. Equivalently, the rows of  $D$  are to be divided into  $K$  disjoint sets (e.g., into red, blue, and green clusters as illustrated here).

(1-2)

Notation:  $A, X, Y$  : matrices or sets (uppercase letters)  
 $\bar{x}, \bar{y}, \bar{\alpha}, \bar{\theta}$  : vectors (lower case letters with a bar)  
 $x, \beta, r, a$  : scalars (lower case letters)

Another ("metric") way to think about clustering is illustrated in 2D.

•  $\bar{x}_i$



One objective function used in clustering is to minimize the sum of all within cluster pairwise distances while also maximizing the sum of distances of each point to the nearest point outside its cluster.

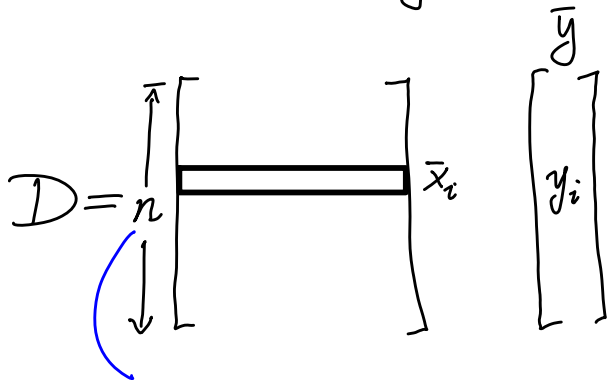
An example: Each  $\bar{x}_i$  represents a consumer, and the entries are the \$ amounts they spent on fruits, meat, toiletries, or other items (in a retail or grocery store). The store may be interested in the clusters to target ads for specific products.

Clustering is considered an unsupervised learning problem, since only the data ( $\bar{x}_i$ ) is used without the membership information (of which cluster each  $\bar{x}_i$  belongs to).

As opposed to classification/regression problems being introduced now, which are supervised learning problems.

## 2. Classification

Here, apart from the  $n \times d$  data matrix  $D$ , we are given also an  $n$ -vector  $\bar{y}$  of labels for each observation.



For instance,  $y_i \in \{-1, 1\}$  (or  $\{0, 1\}$ ) to capture a YES/NO aspect of each observation, e.g., whether customer  $i$  is possibly responsive to ads. In another instance,  $y_i \in \{R, G, B\}$  (3 colors). Note that such labels are categorical, and it is not obvious how to represent these labels using numerical values. The YES/NO aspect could still be modeled as a continuous numerical value coupled with a cut-off value (to determine YES).

The goal is to build a "model" on  $D$  and  $\bar{y}$ , which is the training set, and use that model to predict the  $\bar{y}_t$  values for an external test data set  $D_t$  ( $n_t \times d$ ) for which  $\bar{y}_t$  is not known.

The model can be described as

$$y_i \approx f(\bar{x}_i).$$

The function to be learned,  $f$ , is often parametrized using a weight vector  $\bar{w}$ , and hence we write

$$y_i \approx \bar{w}^T \bar{x}_i$$

(1.4)

For instance, in the case of binary classification with  $y_i \in \{-1, 1\}$ , we could take

$$y_i = \text{sign} \{ f_{\bar{w}}(\bar{x}_i) \}$$

Q. How does one choose  $\bar{w}$ ?

We usually pick  $\bar{w}$  such that the mismatch between  $y_i$  and  $f_{\bar{w}}(\bar{x}_i)$ , i.e., the "loss", is minimized.

The form of  $f$  and this loss function are often carefully chosen/constructed to solve

$$\min_{\bar{w}} \sum_{i=1}^n l(y_i - f_{\bar{w}}(\bar{x}_i))$$

→ loss function

Such classification tasks are considered supervised learning, since we are guided by the  $y_i$ 's (known labels). The goal is to use the built function  $f_{\bar{w}}(\bar{x})$  to predict the labels  $y_i$  for an external or test data set  $D_t$  ( $n_t \times d$ ) (for which the labels are unknown). Since  $D_t$  is not used in the training process, such predictions for test sets are called as generalizations.

The step where we want to minimize a loss function makes classification an optimization problem! Before considering details of optimization, we present a simpler ML problem: regression.

### 3. Regression

We're given  $D$  ( $n \times d$ ) just as in classification, but now,  $y_i \in \mathbb{R}$ , i.e., it is a numerical value (rather than a label). The simplest version of regression is linear regression, where the task is to fit a line through the given set of  $n$  points (in 2D, and a plane in  $d$ -dimensions).

$$y_i = f_{\bar{w}}(\bar{x}_i) = \bar{w}^T \bar{x}_i = \bar{x}_i^T \bar{w}$$

Equivalently,  $\bar{y} = D\bar{w}$ .

And the loss function usually used is the sum of squared errors:

$$J = \frac{1}{2} \|D\bar{w} - \bar{y}\|^2$$

Hence, linear regression becomes the optimization problem:

$$\min_{\bar{w}} \frac{1}{2} \|D\bar{w} - \bar{y}\|^2$$

### Optimization in 1D Calculus

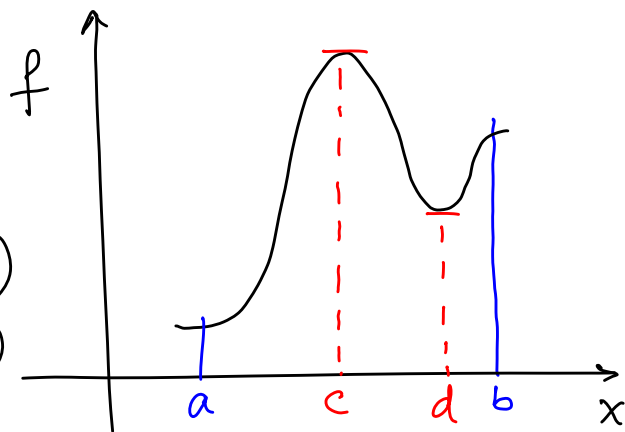
Find  $\min f(x)$  for  $a \leq x \leq b$

\* We set  $f'(x) = 0$  to find critical points.

\*  $f''(x) > 0$ : local minima (e.g.,  $d$ )

$f''(x) < 0$ : local maxima (e.g.,  $c$ )

$f''(x) = 0$ : saddle points



Compare the function values at local minima with those at the end points ( $a$  and  $b$ ) to determine the true ("global") minimum (@  $x = a$  in the figure here).

## Optimization in d-dimensions

1.6

In general, we consider problems of the form

$$\begin{aligned} \min & f(\bar{x}) \\ \text{s.t.} & \bar{g}(\bar{x}) \leq \bar{0} \\ & \bar{h}(\bar{x}) = \bar{0} \end{aligned}$$

Under appropriate assumptions, we can specify optimality conditions (local by default, global when the functions are "nice"), e.g., Karush-Kuhn-Tucker (KKT) conditions.

Corresponding to  $f'(\bar{x})=0$ , we have  $\nabla f = \bar{0}$  for

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}, \text{ the gradient of } f.$$

Back to linear regression:

$$\text{Recall: } \min_{\bar{w}} J = \frac{1}{2} \|D\bar{w} - \bar{y}\|^2$$

$$\Rightarrow \nabla J = D^T D \bar{w} - D^T \bar{y} = \bar{0}$$

$$\Rightarrow \bar{w} = (D^T D)^{-1} D^T \bar{y} \quad \text{closed form expression!}$$

It turns out that this critical point is the (unique) global minimizer of  $J$ . In this sense, linear regression is perhaps the easiest ML problem!

(1.7)

But more generally, we do not get closed form solutions in this form. In fact, we usually use gradient descent to iteratively update  $\bar{w}$ :

We choose an initial  $\bar{w}$  in some way, e.g., randomly. Then, in each iteration, we update

$$\bar{w} \leftarrow \bar{w} - \alpha \nabla J(\bar{w})$$

$\alpha$  is the step size, also called the learning rate. We usually need to choose  $\alpha$  carefully to ensure the iterations converge. We also have to make assumptions about the function  $f$  to guarantee efficient convergence of this gradient-descent method.



# MATH 565: Lecture 2 (01/15/2026)

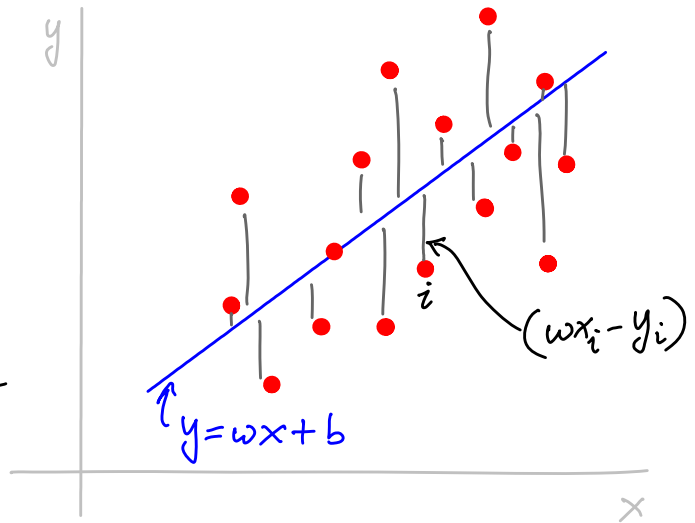
Today: \* regression and extensions  
 \* Regularization  
 \* Support vector machines (SVM)

Recall: Linear regression: given  $D_{n \times d}$ ,  $\bar{y} \in \mathbb{R}^n$ , the goal is to find weight vector  $\bar{w} \in \mathbb{R}^d$  such that  $\bar{y} \approx D\bar{w}$ . In detail, we want to minimize the loss function

$$J = \frac{1}{2} \|D\bar{w} - \bar{y}\|^2$$

$$= \frac{1}{2} \sum_{i=1}^n (y_i - \bar{w}^T \bar{x}_i)^2$$

In 1D, linear regression fits a line  $y = wx + b$  through a given set of  $n$  points  $(x_i, y_i)$  such that the sum of the squared "errors"  $\sum_{i=1}^n (wx_i - y_i)^2$  is minimized.



$J$  is a convex function and hence has a unique minimum. The first order optimality condition (corresponding to  $f'(x) = 0$  in 1D) is given by  $\nabla J = 0$

$$J = \frac{1}{2} (D\bar{w} - \bar{y})^T (D\bar{w} - \bar{y}) = \frac{1}{2} [\bar{w}^T D^T D \bar{w} - 2 \bar{w}^T D^T \bar{y} + \bar{y}^T \bar{y}]$$

$$\text{So, } \nabla J = D^T D \bar{w} - D^T \bar{y} = 0$$

$$\Rightarrow D^T D \bar{w} = D^T \bar{y} \quad \text{--- (1)}$$

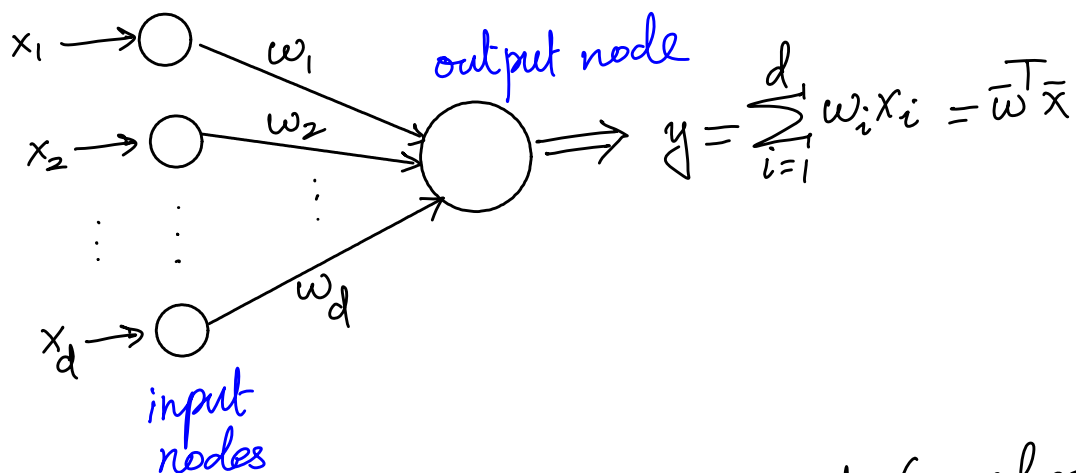
$$\Rightarrow \bar{w} = (D^T D)^{-1} D^T \bar{y} \quad \text{--- (2)}$$



## A Quick aside: Optimization on Computational Graphs

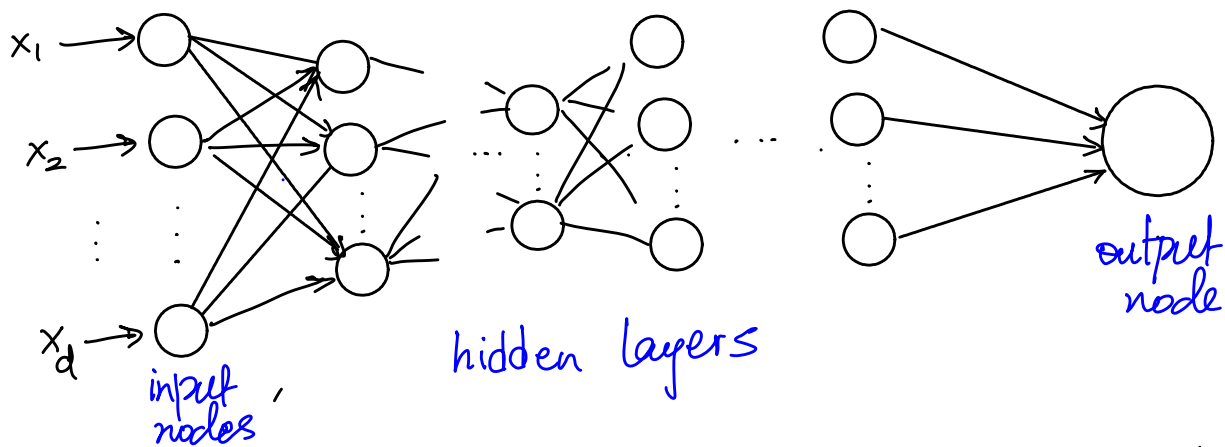
Many machine learning problems can be equivalently posed as problems on computational graphs. Informally, the qualifier "computational" here means the vertices and edges in these graphs are endowed with more general computational steps than, e.g., the weights/costs and edge capacities seen in network flow problems. While we will revisit this topic in detail later on, we give a graph that can be used to model linear regression.

### Linear regression graph model



In general, each node takes in all its inputs (on edges coming in to it) and combines them somehow to send output(s) to all nodes to which it is connected. For linear regression, there is one input node for each  $x_i$  that uses the weight  $w_i$  on its output edge to send  $w_i x_i$  to the single output node. And this output node sums up all its inputs to output  $\sum_i w_i x_i = \bar{w}^T \bar{x}$ .

But the computational graphs could be far more general, with multiple layers of nodes and edges connecting nodes in each layer to many or all nodes in nearby (not necessarily only adjacent) layers.



The functions and data captured on the nodes and edges can also be quite general. As one can imagine, such graphs could capture fairly complicated functions and restrictions. How does one take gradients on computational graphs? We use back propagation. More on this later...

### Back to Regression...

Solving a system of linear equations can be considered as a version of linear regression. If there is a solution, then we get  $J = 0$  (zero loss). But if there is no (exact) solution, we could still get a "best fit" solution (in the least squares sense).

Recall how we found  $\bar{w}$ :

$$\bar{D}^T \bar{D} \bar{w} = \bar{D}^T \bar{y} \quad \text{--- (1)}$$

$$\bar{w} = (\bar{D}^T \bar{D})^{-1} \bar{D}^T \bar{y} \quad \text{--- (2)}$$

How do we compute  $\bar{w}$  efficiently (for large instances)? Computing  $(\bar{D}^T \bar{D})^{-1}$  can be costly, especially when  $n$  and  $d$  are huge (or even large). We usually use QR decomposition of the data matrix  $D$ .

$$D = QR \quad \text{--- (3)}$$

where  $Q$  is  $n \times d$  with orthonormal columns and  $R$  is  $d \times d$  upper triangular.

Hence,  $Q^T Q = I_d$  (identity matrix).

Using (3) in (1) gives

$$R^T \underbrace{Q^T Q}_I R \bar{w} = R^T Q \bar{y}$$

$$\Rightarrow (R^T)^T (R^T R \bar{w} = R^T Q \bar{y}) \quad \text{--- (4)}$$

$$\Rightarrow R \bar{w} = Q^T \bar{y}$$

This is a triangular system and can be solved by back substitution.  
 ↳ as  $R$  is upper triangular

This computation assumes that  $D^T D$  is invertible, which may not always hold — e.g., when  $n < d$  (we do not have enough samples for the given number of dimensions). This is similar to the setting in which a linear system  $A\bar{x} = \bar{b}$  has infinitely many solutions. In this setting, we run the danger of overfitting the (training) data. And when that happens, the trained model may not generalize as well to external test data.

To resolve this situation, we need to do regularization. Intuitively, we need to limit the #  $w_j$ 's that are non-zero...

# Tikhonov Regularization

We modify the loss function to

$$J_w = \frac{1}{2} \|D\bar{w} - \bar{y}\|^2 + \frac{\lambda}{2} \|\bar{w}\|^2$$

for the regularization parameter  $\lambda > 0$ .

Our goal is to somehow enforce several or many of the  $w_j = 0$  (i.e., only a few of them are non-zero). One approach to achieve this goal is to add an explicit constraint of the form

0-norm  $\rightarrow \|\bar{w}\|_0 \leq k$  for  $k \ll d$ .

(# non-zero entries)

But this constraint is hard to enforce, and also destroys the nice structure of  $J$  (convexity). Instead, we add the squared norm penalty  $\frac{\lambda}{2} \|\bar{w}\|^2$  as a regularization term. It achieves the same goal — it forces many  $w_j = 0$  in the optimal solution. But it is also a strongly convex function ( $\lambda > 0$ ), and hence  $J_w$  is so as well ( $\frac{1}{2} \|D\bar{w} - \bar{y}\|^2$  term is convex), implying that  $J_w$  has a unique optimal solution.  $\rightarrow w_j \neq 0$  unless absolutely needed.

Setting  $\nabla J_w = \bar{0}$  for optimality gives

$$\begin{aligned} (D^T D + \lambda I) \bar{w} &= D^T \bar{y} \\ \Rightarrow \bar{w} &= (D^T D + \lambda I)^{-1} D^T \bar{y} \end{aligned}$$

We'll talk about more details of this regularized model later on...

## Modifications of $J_w$

$$\text{In } J_w = \frac{1}{2} \|D\bar{w} - \bar{y}\|_2^2 + \frac{1}{2} \lambda \|\bar{w}\|_2^2, \quad \text{also called "fit"}$$

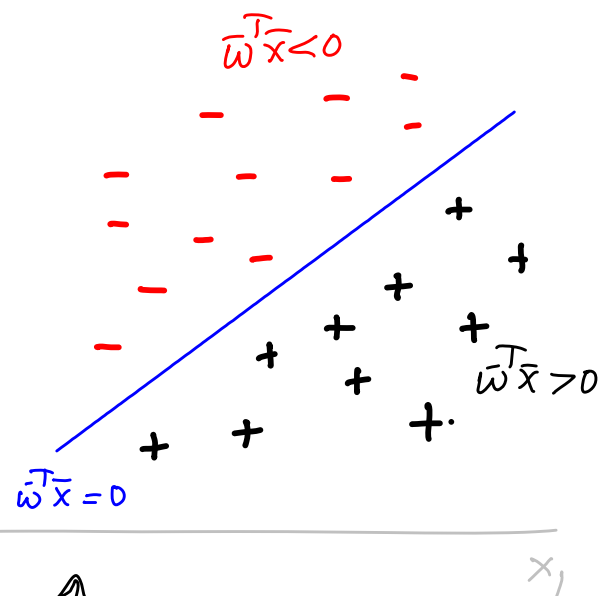
one could consider modifying either (loss/error or regularization) term to an  $L_1$  norm term (instead of  $L_2$ ). There are advantages and disadvantages for using  $L_1$  terms.  $L_1$  norms are piecewise linear, and when they appear in minimization objective functions as here, they can be easily linearized, making computations more efficient. At the same time, the  $L_2$  terms usually have smoother behavior than their  $L_1$  counterparts.

## Binary Classification

We now consider the binary classification problem in a bit more detail in the context of regression. Here, we have  $y_i \in \{-1, 1\}$  and the goal is to "separate" the  $+1$  instances "as best as possible" from the  $-1$  instances.

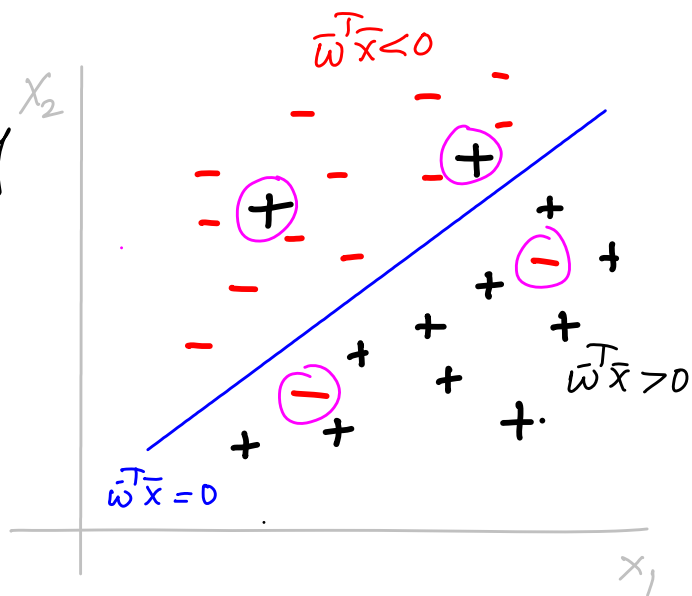
A direct generalization of linear regression is to find  $\bar{w}$  s.t.  $\bar{w}^T \bar{x} = 0$  is the separating line (in 2D) or hyperplane (in  $d$ -dim), with the  $\bar{w}^T \bar{x} > 0$  side capturing all the  $+1$  instances and the other side containing all the  $-1$  instances.

This may work well in well-separated, i.e., easy, instances as illustrated here.



But when the  $+1$  and  $-1$  instances are not well-separated (as seen here), we would want to modify  $J_w$  (loss function) to capture the violations of good separations. Here is one approach. We write

$$J_w = \frac{1}{2} \sum_{i=1}^n (y_i - \bar{w}^T \bar{x}_i)^2 + \frac{1}{2} \|\bar{w}\|^2$$



Using the fact that  $y_i^2 = 1$  (as  $y_i \in \{-1, 1\}$ ), we write

$$J_w = \frac{1}{2} \sum_{i=1}^n y_i^2 (y_i - \bar{w}^T \bar{x}_i)^2 + \frac{1}{2} \|\bar{w}\|^2$$

$$= \frac{1}{2} \sum_{i=1}^n (y_i^2 - y_i \bar{w}^T \bar{x}_i)^2 + \frac{1}{2} \|\bar{w}\|^2$$

$$= \frac{1}{2} \sum_{i=1}^n (1 - y_i \bar{w}^T \bar{x}_i)^2 + \frac{1}{2} \|\bar{w}\|^2$$

(again, using  $y_i^2 = 1$ ).

The first term will penalize violations of good separation. At the same time, it will also penalize cases that are "extremely" well-separated! For instance, if  $y_i = +1$  and  $\bar{w}^T \bar{x}_i = 100$ , which is a very good separation, the term in  $J_w$  will be  $(1 - 1 \cdot 100)^2 = 99^2 = 9801$ , which is huge!

Hence, we want to set the contribution of well-separated instances (to  $J_w$ ) as zero — this is the idea used in support vector machines (SVMs)!

# Support Vector Machines (SVM)

(2-7)

We set

$$J_{L_2\text{-svm}} = \frac{1}{2} \sum_{i=1}^n (\max\{0, [1 - y_i(\bar{w}^T \bar{x}_i)]\})^2 + \frac{1}{2} \|\bar{w}\|^2$$

A (computationally) better option is to consider an  $L_1$ -SVM loss.

We rewrite this  $L_1$ -SVM problem as a slightly modified optimization problem, as presented below. For many people working in optimization, this SVM model is a natural way to start exploring machine learning problems, as it is a convex optimization problem.

$$\begin{aligned} \min_{\bar{w}, \bar{\epsilon}, b} \quad & J_{\text{svm}} = C \sum_{i=1}^n \bar{\epsilon}_i + \frac{1}{2} \|\bar{w}\|^2 \\ \text{s.t.} \quad & y_i (\bar{w}^T \bar{x}_i + b) \geq 1 - \bar{\epsilon}_i, \quad i=1, \dots, n \\ & \bar{\epsilon}_i \geq 0 \end{aligned}$$

We will check the details of this model in the next lecture...



# MATH 565: Lecture 3 (01/20/2026)

Today: \* more on SVM  
\* Taylor expansion  
\* local optimality conditions

Recall SVM  $\min_{\bar{w}, b, \bar{\epsilon}} J = C \sum_{i=1}^n \epsilon_i + \frac{1}{2} \|\bar{w}\|^2 \quad (C > 0)$

s.t.  $y_i (\bar{w}^T \bar{x}_i + b) \geq 1 - \epsilon_i, \quad i=1, \dots, n$   
 $\epsilon_i \geq 0 \quad \forall i$

Let's examine the unified (main) constraints in detail:

When  $y_i = +1$ , we get  $\bar{w}^T \bar{x}_i + b \geq 1 - \epsilon_i$

e.g., if  $\bar{w}^T \bar{x}_i + b = 0.7$ , then  $\epsilon_i \geq 0.3$

But if  $\bar{w}^T \bar{x}_i + b = 2$ , then  $\epsilon_i = 0$  works

$\epsilon_i$ : measures by how much the  $i$ th sample violates well-separatedness

When  $y_i = -1$ , we get  $\bar{w}^T \bar{x}_i + b \leq -1 + \epsilon_i$ .

e.g., if  $\bar{w}^T \bar{x}_i + b = -3$ ,  $\epsilon_i = 0$  in the opt. soln.

But if  $\bar{w}^T \bar{x}_i + b = 0.5$ , we need  $\epsilon_i \geq 1.5$

Can also regularize  $b$  (intercept term).

Or, take  $\bar{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$  as (d+1)-variable vector, and

write  $J = \min_{\bar{w}, \bar{\epsilon}} C \sum_{i=1}^n \epsilon_i + \frac{1}{2} \|\bar{w}\|^2$

s.t.  $y_i (\bar{w}_{1:d}^T \bar{x}_i + w_0) \geq 1 - \epsilon_i, \quad i=1, \dots, n$   
 $\epsilon_i \geq 0 \quad \forall i.$

# Taylor Expansion

3.2

In 1D, Taylor expansion of  $f(x)$  at  $x=a$  is

$$f(x) = f(a) + (x-a)f'(a) + \frac{(x-a)^2}{2!}f''(a) + \dots + \frac{(x-a)^r}{r!} \left. \frac{d^r f(x)}{dx^r} \right|_{x=a} + \dots$$

When  $|x-a|$  is small, we could take

$$f(x) \approx f(a) + (x-a)f'(a) + \frac{(x-a)^2}{2!}f''(a)$$

In d-dimensions, Taylor expansion of  $f(\bar{x})$  at  $\bar{x}=\bar{a}$

$$f(\bar{x}) = f(\bar{a}) + \sum_{i=1}^d (x_i - a_i) \left[ \frac{\partial f}{\partial x_i} \right] \Big|_{\bar{x}=\bar{a}} + \sum_i \sum_j \frac{(x_i - a_i)(x_j - a_j)}{2!} \left[ \frac{\partial^2 f}{\partial x_i \partial x_j} \right] \Big|_{\bar{x}=\bar{a}} + \dots$$

Equivalently,

$$f(\bar{x}) = f(\bar{a}) + [\bar{x} - \bar{a}]^T \nabla f(\bar{a}) + [\bar{x} - \bar{a}]^T H f(\bar{a}) [\bar{x} - \bar{a}] + \dots$$

# Local Optimality Conditions (in 1D)

(3.3)

**Lemma 1**  $f(x)$  is a minimum value at  $x=x_0$  with respect to its immediate locality if  $\xrightarrow{\text{sufficient condition}}$   
 $\underline{f'(x_0)=0}$  and  $\underline{f''(x_0) > 0}$ .  
 $\swarrow$  first order condition for optimality  $\searrow$  second order condition for optimality.

Proof Consider Taylor expansion of  $f(\cdot)$  at  $x_0$  ( $x=x_0+\Delta$ )  
 $f(x_0+\Delta) \approx f(x_0) + \Delta \underbrace{f'(x_0)}_{=0} + \frac{\Delta^2}{2} \underbrace{f''(x_0)}_{>0}$  for  $|\Delta| \ll 1$ .  
Small enough

$\Rightarrow$  Under first and second order optimality conditions,  
 $f(x_0+\Delta) > f(x_0)$  for  $|\Delta|$  small enough.  $\square$

The first order condition ( $f'(x)=0$ ) is typically solved using gradient descent.

Step 0. Start  $x=x_0$  (randomly chosen)

Step k.  $x_k \leftarrow x_{k-1} - \alpha f'(x_{k-1})$   
(in general,  $x \leftarrow x - \alpha f'(x)$ )

$\alpha$ : learning rate (step size)

Changing  $x$  by  $\Delta x = -\alpha f'(x)$ .

We're changing  $x$  along the "steepest descent" direction (trivial in 1D, nontrivial in d-dim.).

This change will reduce  $f(x)$  for 'small' values of  $\alpha$ :

$$\begin{aligned} f(x+\delta x) &\approx f(x) + \delta x f'(x) \\ &= f(x) - \alpha [f'(x)]^2 \\ &< f(x) \end{aligned}$$

---

Example

$$f(x) = x^2 \sin(x) + x$$

$$f'(x) = 2x \sin(x) + x^2 \cos(x) + 1 \stackrel{?}{=} 0$$

$$f''(x) = (2-x^2) \sin(x) + 4x \cos(x)$$

## Local Optimality in d-dim

Variables:  $\bar{w}$  ( $w_1, \dots, w_d, w_b$ ),  $\xi_i, \dots$

Obj. fn: loss function  $J$  e.g.,  $J = \frac{1}{2} \|D\bar{w} - \bar{y}\|^2 + \frac{\lambda}{2} \|\bar{w}\|^2$

1<sup>st</sup> order condition:  $\nabla J = \bar{0}$   $\begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_d} \end{bmatrix} = \bar{0}$

2<sup>nd</sup> order condition:  $H \succ 0$  Hessian is positive definite  
i.e.,  $\bar{x}^T H \bar{x} > 0 \quad \forall \bar{x} \in \mathbb{R}^d / \{\bar{0}\}$

Taylor expansion:  $\epsilon > 0$

$$J(\bar{w}_0 + \epsilon \bar{v}) \approx J(\bar{w}_0) + \epsilon \bar{v}^T \underbrace{\nabla J(\bar{w}_0)}_{=0} + \frac{\epsilon^2}{2} \underbrace{\bar{v}^T H \bar{v}}_{>0 \quad \forall \bar{v}}$$