# COMP 551 – Applied Machine Learning
## Lecture 2: Linear regression

**Instructor**:  Joelle Pineau (*jpineau@cs.mcgill.ca)*

**Class web page**: *www.cs.mcgill.ca/~jpineau/comp551*

# Today's Quiz (informal)

Write down the 3 most useful insights you gathered from the article:

*"A Few Useful Things to Know About Machine Learning"*.

- domain-specific ≡ feature selection
- no free lunch.
- simplicity ≠ usefulness

# Supervised learning

- Given a set of **training examples**: $x_i = <x_{i1}, x_{i2}, x_{i3}, \ldots, x_{in}, y_i>$

  $x_{ij}$ is the $j^{th}$ feature of the $i^{th}$ example

  $y_i$ is the desired **output** (or **target)** for the $i^{th}$ example.

  $X_j$ denotes the $j^{th}$ feature.

- We want to learn a function $f : X_1 \times X_2 \times \ldots \times X_n \rightarrow Y$
  which maps the input variables onto the output domain.

| tumor size | texture | perimeter | . . . | outcome | time |
|------------|---------|-----------|-------|---------|------|
| 18.02 | 27.6 | 117.5 | | N | 31 |
| 17.99 | 10.38 | 122.8 | | N | 61 |
| 20.29 | 14.34 | 135.1 | | R | 27 |
| . . . | | | | | |

# Supervised learning

- Given a dataset $X \times Y$, find a function: $f : X \rightarrow Y$ such that $f(\boldsymbol{x})$ is a good predictor for the value of $y$.

- Formally, $f$ is called the **hypothesis**.

- Output $Y$ can have many types:

  – If $Y = \Re$, this problem is called **regression**.

  – If $Y$ is a finite discrete set, the problem is called **classification**.

  – If $Y$ has 2 elements, the problem is called **binary classification**.

# Prediction problems

- The problem of predicting <u>tumour recurrence</u> is called:

**<u>classification</u>**

- The problem of predicting the <u>time of recurrence</u> is called:

**<u>regression</u>**

- Treat them as two separate supervised learning problems.

| tumor size | texture | perimeter | . . . | outcome | time |
|------------|---------|-----------|-------|---------|------|
| 18.02 | 27.6 | 117.5 | | N | 31 |
| 17.99 | 10.38 | 122.8 | | N | 61 |
| 20.29 | 14.34 | 135.1 | | R | 27 |
| . . . | | | | | |

# Variable types

- **Quantitative**, often real number measurements.

  – Assumes that similar measurements are similar in nature.

- **Qualitative**, from a set (categorical, discrete).

  – E.g. {Spam, Not-spam}

- **Ordinal**, also from a discrete set, without metric relation, but that

  allows ranking.

  – E.g. {first, second, third}

# The i.i.d. assumption

- In supervised learning, the examples $x_i$ in the training set are assumed to be independently and identically distributed.

  - Independently: Every $x_i$ is freshly sampled according to some probability distribution $D$ over the data domain $X$.

  - Identically: The distribution $D$ is the same for all examples.

- Why?

# Empirical risk minimization

For a given function class $F$ and training sample $S$,

- Define a notion of error (*left intentionally vague for now*):

$$L_S(f) = \text{\# mistakes made on the sample } S$$

- Define the Empirical Risk Minimization (ERM):

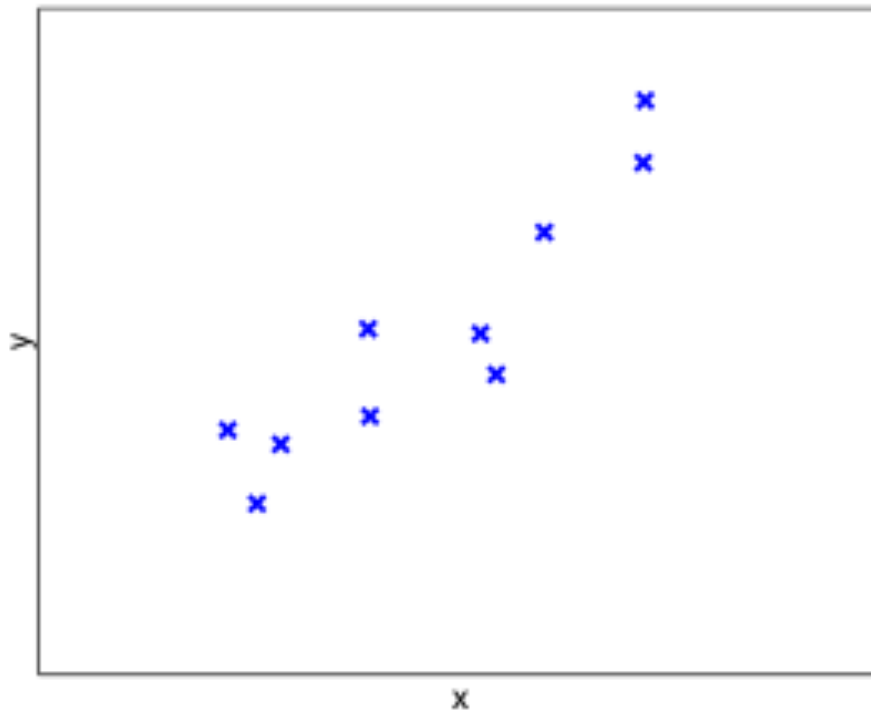$$ERM_F(S) = argmin_{f \text{ in } F} L_S(f)$$

where *argmin* returns the function $f$ (or set of functions) that achieves the minimum loss on the training sample.

- Easier to minimize the error with i.i.d. assumption.

# A regression problem

- What <u>hypothesis class</u> should we pick?



| *Observe* | *Predict* |
|:---:|:---:|
| *x* | *y* |
| 0.86 | 2.49 |
| 0.09 | 0.83 |
| -0.85 | -0.25 |
| 0.87 | 3.10 |
| -0.44 | 0.87 |
| -0.43 | 0.02 |
| -1.1 | -0.12 |
| 0.40 | 1.81 |
| -0.96 | -0.83 |
| 0.17 | 0.43 |

# Linear hypothesis

- Suppose $Y$ is a **linear function** of $X$:

$$f_{\boldsymbol{W}}(\boldsymbol{X}) \quad = \quad w_0 + w_1 x_1 + \ldots + w_m x_m$$

$$= \quad w_0 + \sum_{j=1:m} w_j x_j$$

- The $w_j$ are called **parameters** or **weights**.

- To simplify notation, we add an attribute $x_0=1$ to the $m$ other attributes (also called **bias term** or **intercept**).

**How should we pick the *weights*?**

# Least-squares solution method

- The linear regression problem: $f_{\boldsymbol{w}}(X) = w_0 + \sum_{j=1:m} w_j x_j$

  where $m$ = the dimension of observation space, i.e. number of features.

$$f_w(X) = w_0 + \sum_{j=1:m} w_j x_j$$

- **Goal**: Find the **best** linear model given the data.

- Many different possible **evaluation** criteria!

- Most common choice is to find the $\boldsymbol{w}$ that minimizes:
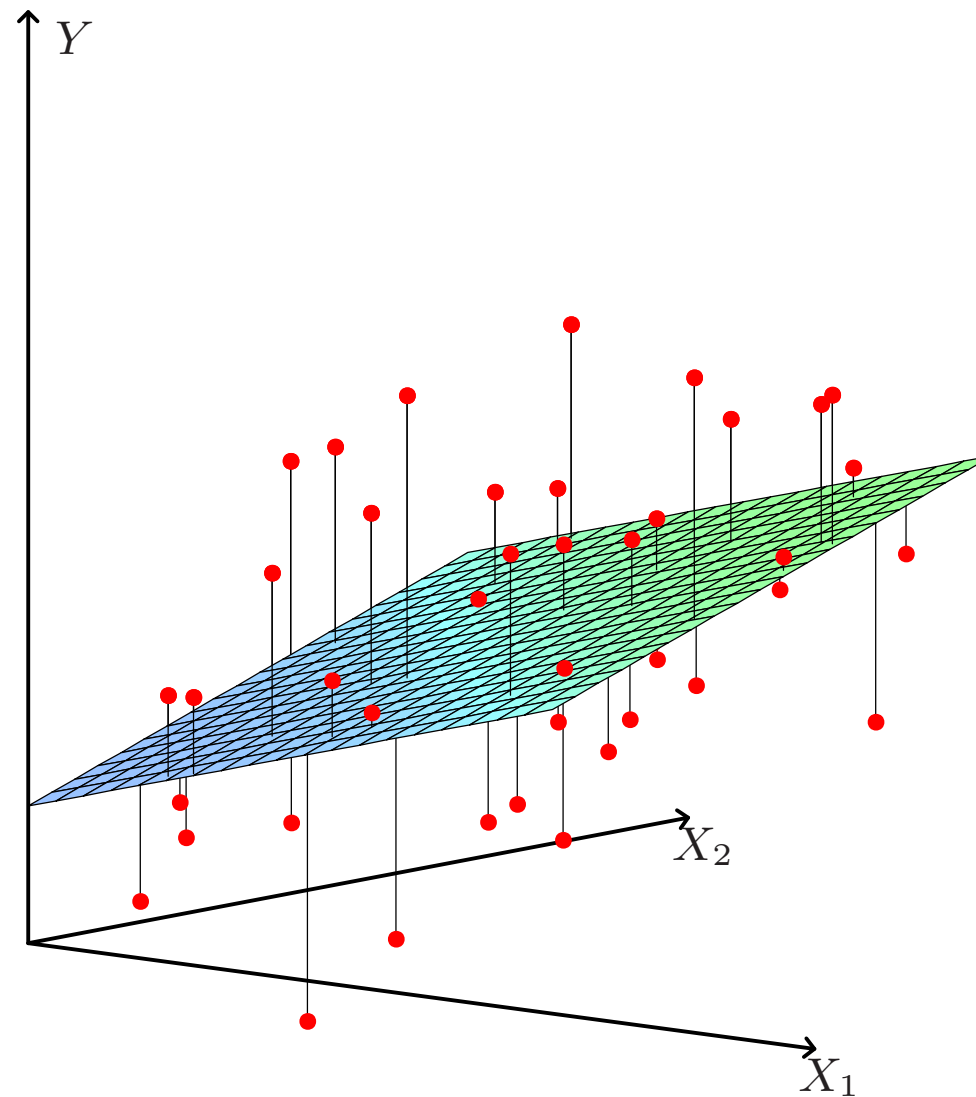
$$Err(w) = \sum_{i=1:n} ( y_i - \boldsymbol{w}^T \boldsymbol{x}_i )^2$$

  (A note on notation: Here $\boldsymbol{w}$ and $\boldsymbol{x}$ are column vectors of size $m+1$.)

$$Err(w) = \sum_{i=1:n} ( y_i - w^T x_j )^2$$

# Least-squares solution for $X \in \mathcal{R}^2$

# Least-squares solution method

- Re-write in matrix notation: $f_{\boldsymbol{w}}(X) = X\boldsymbol{w}$

$$Err(\boldsymbol{w}) = (Y - X\boldsymbol{w})^T(Y - X\boldsymbol{w})$$

where   $X$ is the $n \times m$ matrix of input data,

$Y$ is the $n \times 1$ vector of output data,

$\boldsymbol{w}$ is the $m \times 1$ vector of weights.

- To minimize, take the derivative w.r.t. $\boldsymbol{w}$:

$$\partial Err(\boldsymbol{w})/\partial \boldsymbol{w} = -2\, X^T (Y - X\boldsymbol{w})$$

  – You get a system of $m$ equations with $m$ unknowns.

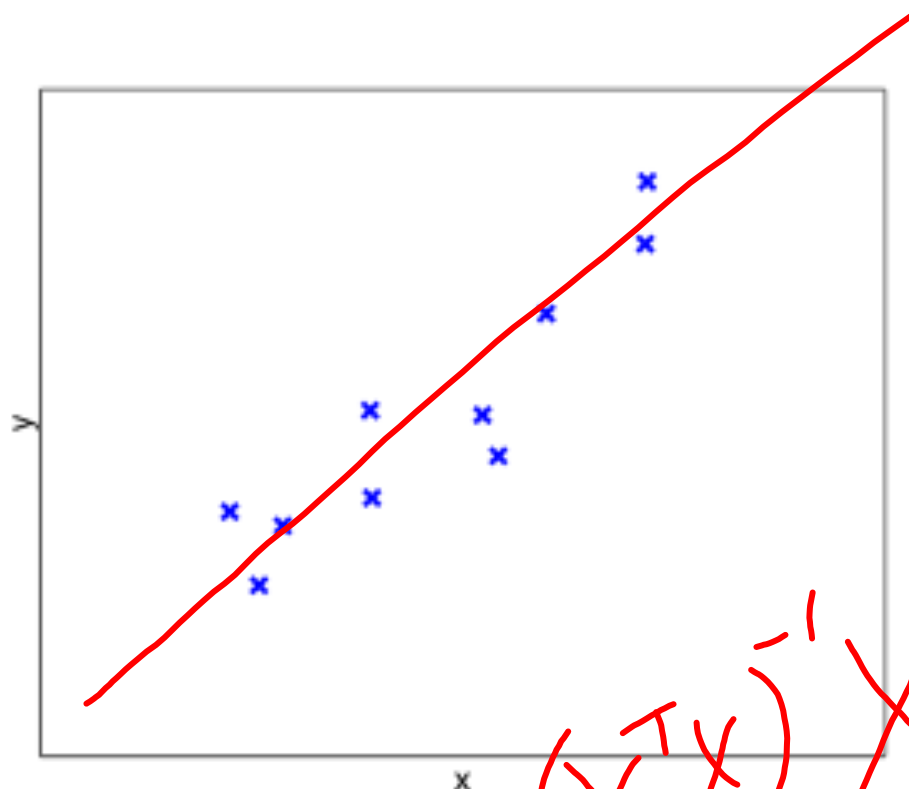- Set these equations to 0:    $X^T (Y - X\boldsymbol{w}) = 0$

# Least-squares solution method

- We want to solve for **w**:      $X^T(Y - X\mathbf{w}) = 0$

- Try a little algebra:      $X^T Y = X^T X \mathbf{w}$

  $$\hat{\mathbf{w}} = (X^T X)^{-1} X^T Y$$

  ($\hat{\mathbf{w}}$ denotes the estimated weights)

$$\hat{w} = (X^T X)^{-1} X^T Y$$

- The fitted data:      $\hat{Y} = X\hat{\mathbf{w}} = X(X^T X)^{-1} X^T Y$

- To predict new data $X' \rightarrow Y'$:      $Y' = X'\hat{\mathbf{w}} = X'(X^T X)^{-1} X^T Y$

| $x$ | $y$ |
|---|---|
| 0.86 | 2.49 |
| 0.09 | 0.83 |
| -0.85 | -0.25 |
| 0.87 | 3.10 |
| -0.44 | 0.87 |
| -0.43 | 0.02 |
| -1.10 | -0.12 |
| 0.40 | 1.81 |
| -0.96 | -0.83 |
| 0.17 | 0.43 |

$$w = (X^T X)^{-1} X Y$$

What is a plausible estimate of **w** ?     **Try it!**
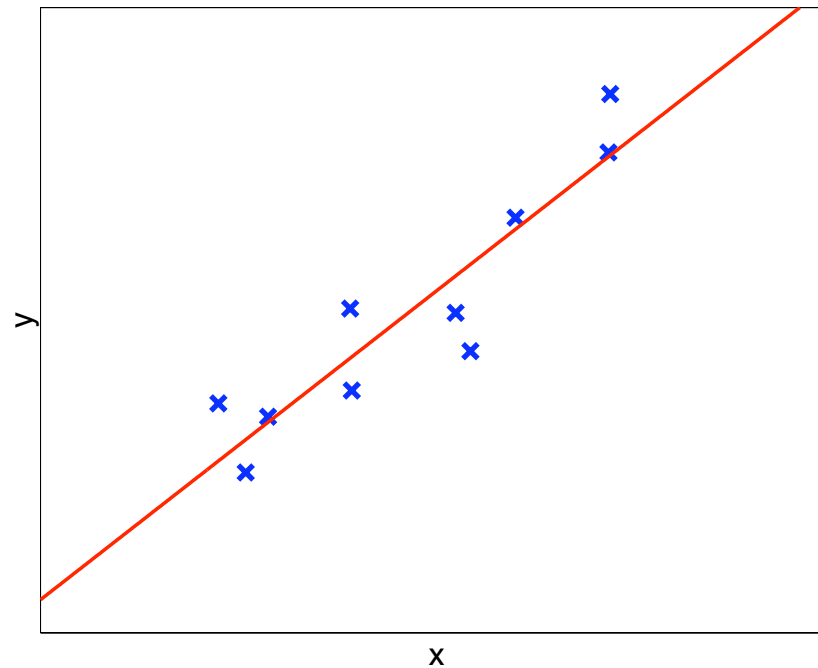
# Data matrices

$$X^T X =$$

$$\begin{bmatrix} 0.86 & 0.09 & -0.85 & 0.87 & -0.44 & -0.43 & -1.10 & 0.40 & -0.96 & 0.17 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 0.86 & 1 \\ 0.09 & 1 \\ -0.85 & 1 \\ 0.87 & 1 \\ -0.44 & 1 \\ -0.43 & 1 \\ -1.10 & 1 \\ 0.40 & 1 \\ -0.96 & 1 \\ 0.17 & 1 \end{bmatrix}$$

$$n \times m \quad \sim 2 \times 10$$

$$= \begin{bmatrix} 4.95 & -1.39 \\ -1.39 & 10 \end{bmatrix}$$

$$m^2$$

# Data matrices

$$X^T Y =$$

$$
\begin{bmatrix}
0.86 & 0.09 & -0.85 & 0.87 & -0.44 & -0.43 & -1.10 & 0.40 & -0.96 & 0.17 \\
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1
\end{bmatrix}
\times
\begin{bmatrix}
2.49 \\
0.83 \\
-0.25 \\
3.10 \\
0.87 \\
0.02 \\
-0.12 \\
1.81 \\
-0.83 \\
0.43
\end{bmatrix}
$$

$$
=
\begin{bmatrix}
6.49 \\
8.34
\end{bmatrix}
$$

# Solving the problem

$$\mathbf{w} = (X^T X)^{-1} X^T Y = \begin{bmatrix} 4.95 & -1.39 \\ -1.39 & 10 \end{bmatrix}^{-1} \begin{bmatrix} 6.49 \\ 8.34 \end{bmatrix} = \begin{bmatrix} 1.60 \\ 1.05 \end{bmatrix}$$

$\beta_1$

$\beta_0$

So the best fit line is $y = 1.60x + 1.05$.

# Interpreting the solution

- Linear fit for a prostate cancer dataset

  – Features $X$ = {lcavol , lweight, age, lbph, svi, lcp, gleason, pgg45}

  – Output $y$ = level of PSA (an enzyme which is elevated with cancer).

  – High coefficient weight (in absolute value) = important for prediction.

| Term | Coefficient | Std. Error |
|---|---|---|
| Intercept | $w_0$ = 2.46 | 0.09 |
| lcavol | 0.68 | 0.13 |
| lweight | 0.26 | 0.10 |
| age | $-0.14$ | 0.10 |
| lbph | 0.21 | 0.10 |
| svi | 0.31 | 0.12 |
| lcp | $-0.29$ | 0.15 |
| gleason | $-0.02$ | 0.15 |
| pgg45 | 0.27 | 0.15 |

# Computational cost of linear regression

$$= nm^2 + m^3$$

- What operations are necessary?

  - Overall: 1 matrix inversion + 3 matrix multiplications

  - $X^T X$          (other matrix multiplications require fewer operations.)

    - $X^T$ is *mxn* and $X$ is *nxm*, so we need $nm^2$ operations.

  - $(X^T X)^{-1}$

    - $X^T X$ is *mxm*, so we need $m^3$ operations.

- We can do linear regression in polynomial time, but handling large datasets (many examples, many features) can be problematic.

*Gradient Solution*

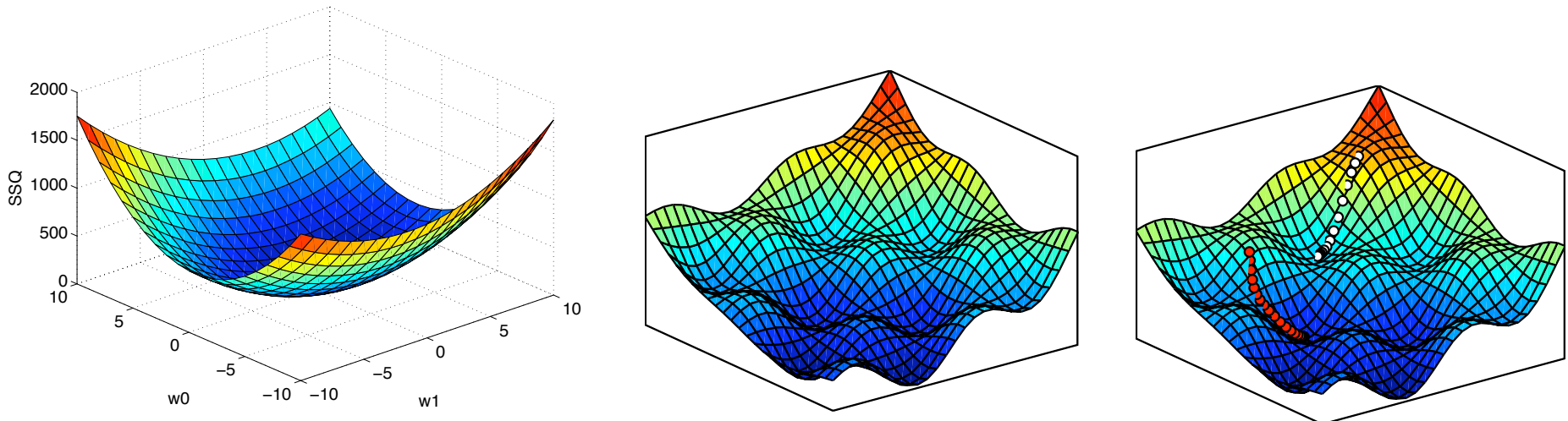# An alternative for minimizing mean-squared error (MSE)

- Recall the least-square solution: $\hat{w} = (X^TX)^{-1}X^T Y$

- What if $X$ is too big to compute this explicitly (e.g. $m \sim 10^6$)?

- Go back to the gradient step: $Err(w) = (Y - Xw)^T(Y - Xw)$

$$\partial Err(w)/\partial w = -2 X^T (Y - Xw)$$

$$\partial Err(w)/\partial w = 2(X^TXw - X^TY)$$

# Gradient-descent solution for MSE

- Consider the error function:



- The gradient of the error is a vector indicating the direction to the minimum point.

- Instead of directly finding that minimum (using the closed-form equation), we can take small steps towards the minimum.

**Gradient descent algorithm**

# Gradient-descent solution for MSE

- We want to produce a sequence of weight solutions, $w_0, w_1, w_2 \dots$,

  such that: $Err(w_0) > Err(w_1) > Err(w_2) > \dots$

- The algorithm:

  > Given an initial weight vector $w_0$,
  >
  > Do for $k=1, 2, \dots$
  >
  > $$w_{k+1} = w_k - \alpha_k\, \partial Err(w_k)/\partial w_k$$
  >
  > End when $|w_{k+1} - w_k| < \varepsilon$

- Parameter $\alpha_k > 0$ is the step-size (or <u>learning rate</u>) for iteration $k$.

# Convergence

- Convergence depends in part on the $\alpha_k$.

- If steps are too **large**: the $w_k$ may oscillate forever.
    - This suggests that $\alpha_k \to 0$ as $k \to \infty$.

- If steps are too **small**: the $w_k$ may not move far enough to reach a local minimum.

# Robbins-Monroe conditions

- The $\alpha_k$ are a Robbins-Monroe sequence if:

$$\sum_{k=0:\infty} \alpha_k = \infty$$

$$\sum_{k=0:\infty} \alpha_k^2 < \infty$$

- These conditions are sufficient to ensure convergence of the $w_k$ to a **local minimum** of the error function.

    E.g. $\alpha_k = 1 / (k + 1)$                 (averaging)

    E.g. $\alpha_k = 1/2$ for $k = 1, ..., T$

            $\alpha_k = 1/2^2$ for $k = T+1, ..., (T+1)+2T$
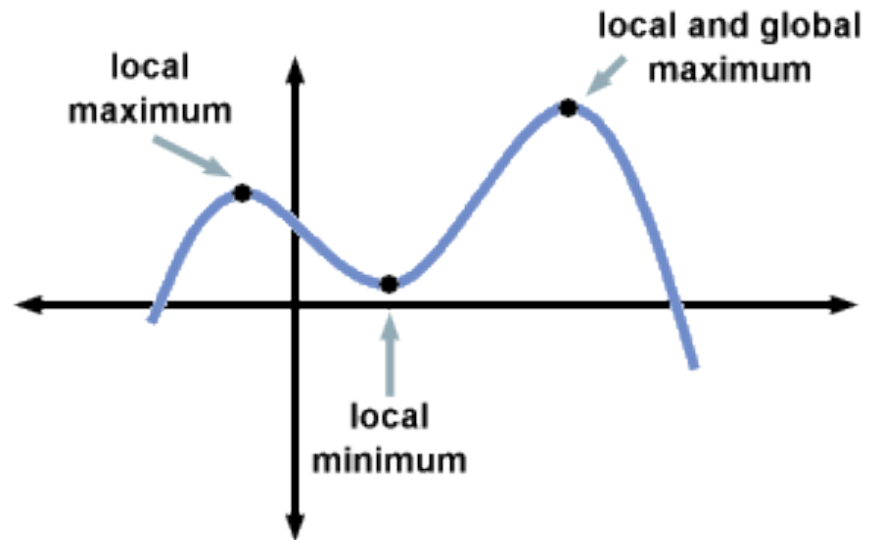
            etc.

# Local minima

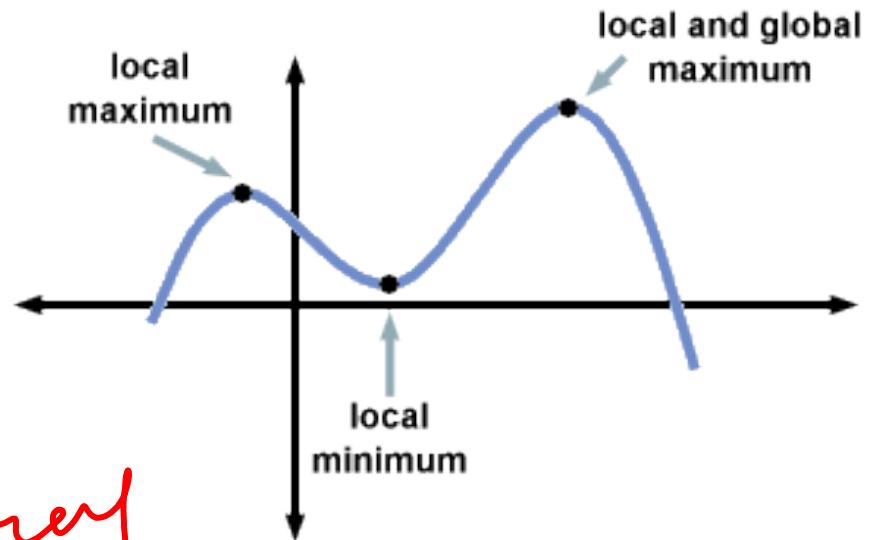- Convergence is **NOT** to a global minimum, only to local minimum.



- The blue line represents the error function. There is *no guarantee* regarding the amount of error of the weight vector found by gradient descent, compared to the globally optimal solution.

# Local minima

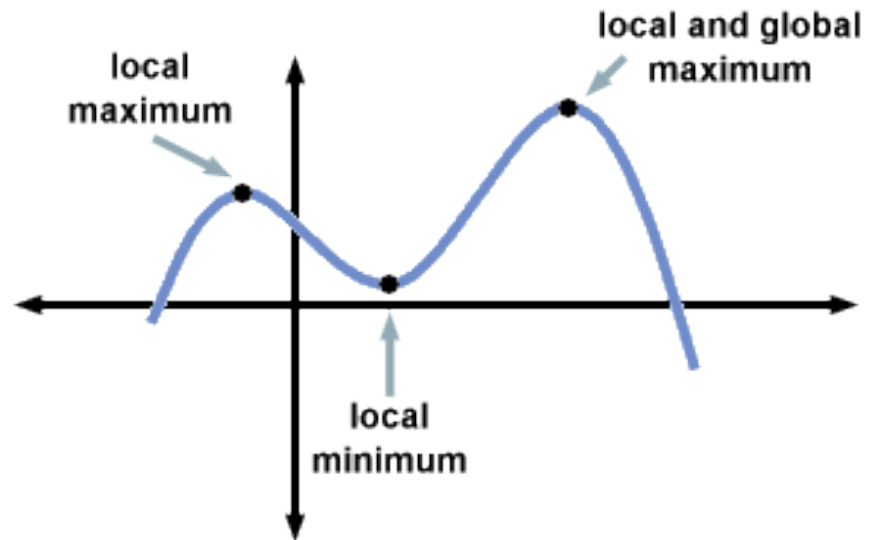- Convergence is **NOT** to a global minimum, only to local minimum.



*linear has no local*

- For linear function approximations using Least-Mean Squares (LMS) error, this is not an issue: only **ONE** global minimum!

  – Local minima affects many other function approximators.

# Local minima

- Convergence is **<u>NOT</u>** to a global minimum, only to local minimum.



- For linear function approximations using Least-Mean Squares (LMS) error, this is not an issue: only **ONE** global minimum!

  – Local minima affects many other function approximators.

- *Repeated random restarts* can help (in all cases of gradient search).

# A 3rd optimization method: QR decomposition (optional)

- Consider the usual criteria: $X^T(Y - Xw)=0$

- Assume $X$ can be decomposed: $X = QR$

  where $Q$ is an $n \times m$ <u>orthogonal</u> matrix (i.e. $Q^TQ=I$), and $R$ is an $m \times m$ upper triangular matrix.

- Replace $X$ in equation above: $(QR)^TY = (QR)^T(QR)w$

- Distribute the transpose: $R^TQ^TY = R^TQ^TQRw$

- Let $Q^TQ=I$ and multiply by $(R^T)^{-1}$ $\quad Q^TY = Rw$

- Solution: $\hat{w} = R^{-1}Q^TY$ $\qquad$ The fitted outputs are: $\hat{Y} = QQ^TY$

- This method is more numerically stable than others, and $R^{-1}$ is fast to compute because upper triangular.

- Alternately, we can use **singular value decomposition**.

# What you should know

- Definition and characteristics of a supervised learning problem.

- Linear regression (hypothesis class, cost function, algorithm).

- Closed-form least-squares solution method (algorithm, computational complexity, stability issues).

- Gradient descent method (algorithm, properties).

# To-do

- Reproduce the linear regression example (slides 15-18), solving it using the software of your choice.

- Suggested complementary readings:

  - Ch.2 (Sec. 2.1-2.4, 2.9) of Hastie et al.

  - Ch.3 of Bishop.

  - Ch.9 of Shalev-Schwartz et al.

- Write down **midterm** date in agenda:  Nov. 22, 6-8pm, Leacock 132.

- Tutorial times (appearing soon): www.cs.mcgill.ca/~jpineau/comp551/schedule.html

- Office hours (confirmed): www.cs.mcgill.ca/~jpineau/comp551/syllabus.html