**Lecture 06:**

**Optimization**
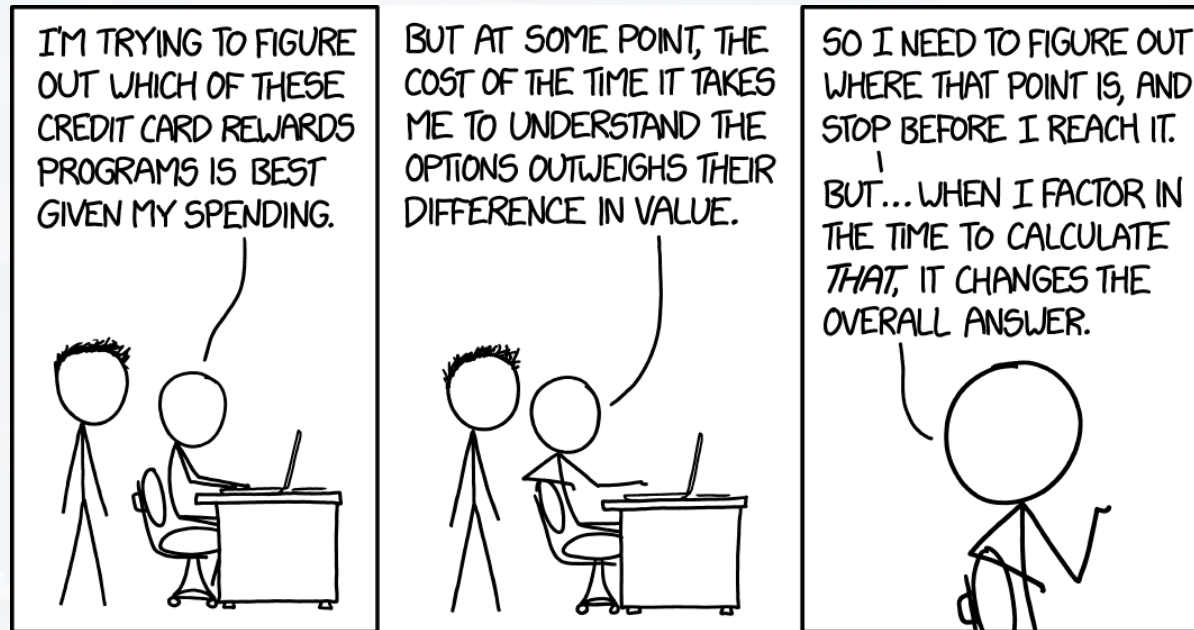
Markus Hohle

University California, Berkeley

**Machine Learning Algorithms**

MSSE 277B, 3 Units

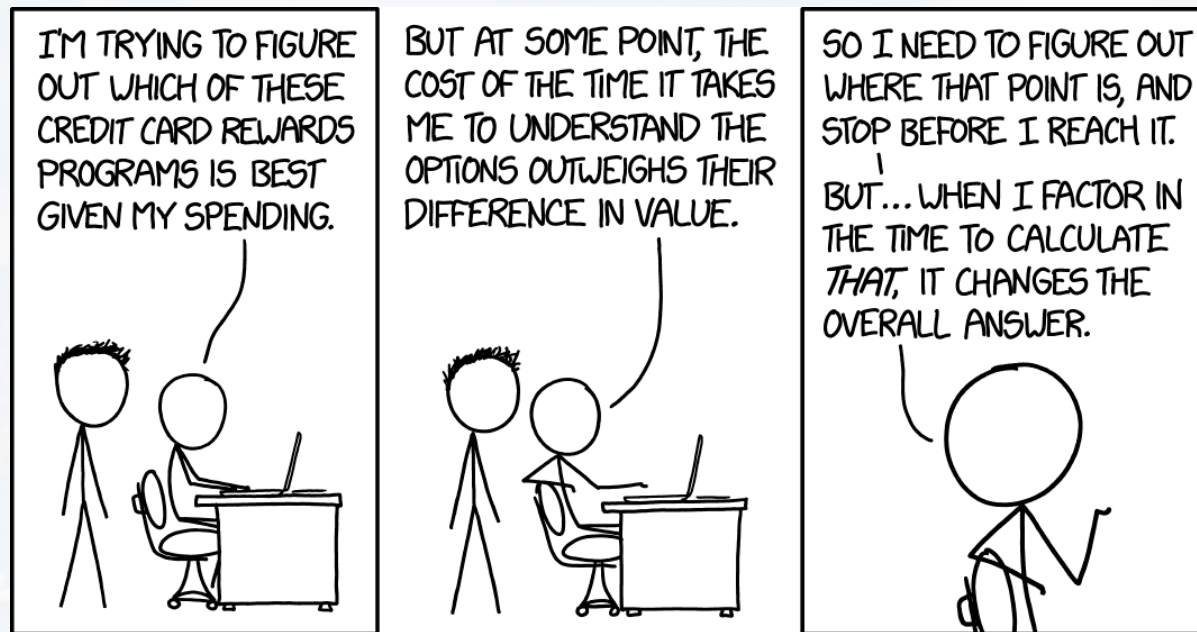Fall 2024

<u>Outline</u>

**- The Problem**

**- Gradient Descent**

- Vanilla
- Learning Rate Schedule
- Momentum
- L1 and L2
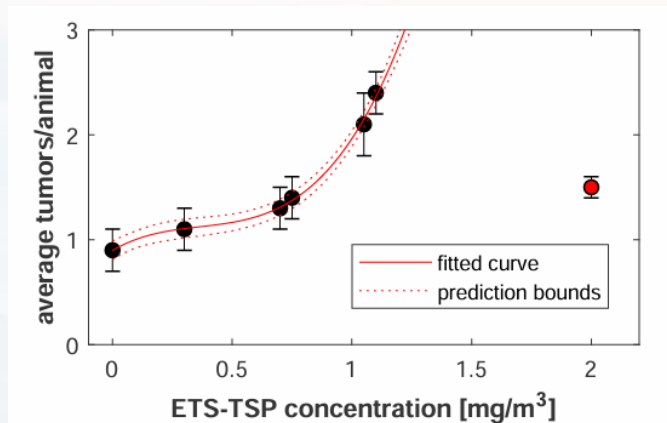- More Finetuning

Outline

**- The Problem**

- Gradient Descent

    - Vanilla
    - Learning Rate Schedule
    - Momentum
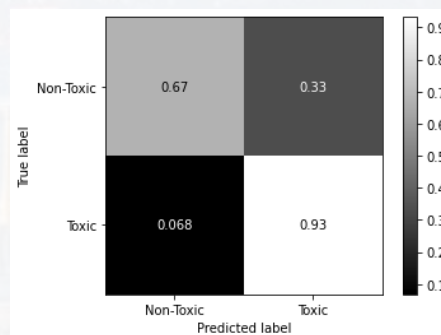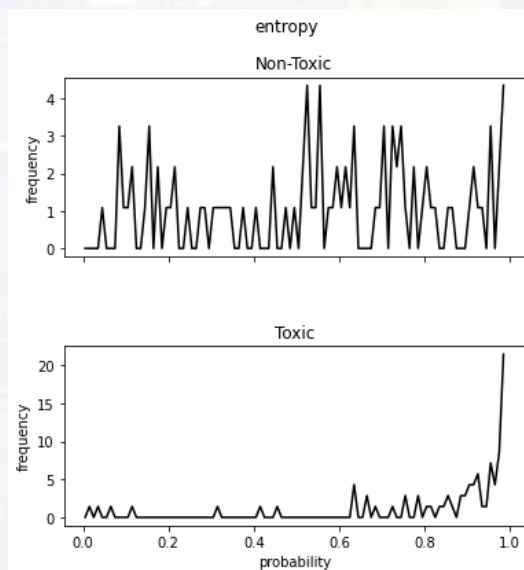    - L1 and L2
    - More Finetuning

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

**regression, e. g. curve fitting**



minimize:     $\chi_{red}^2 = \dfrac{1}{N-p-1} \displaystyle\sum_{i=1}^{N} \dfrac{(\bar{y}(model)_i - y_i)^2}{\sigma_i^2}$

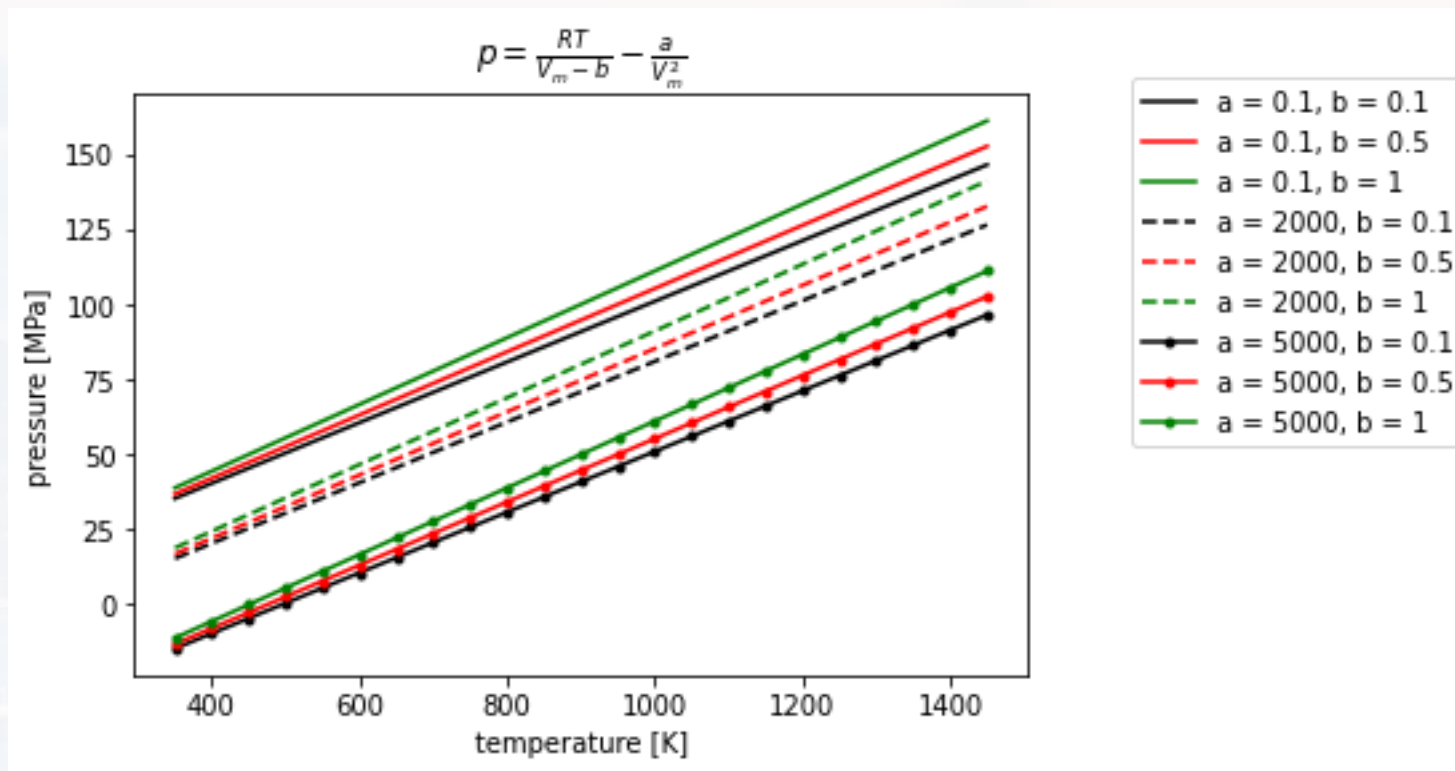**classification**



maximize: accuracy

minimize: cross entropy     $S = -\displaystyle\sum_i p(true)_i \cdot \ln p(model)_i$

# **Optimization**:

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)



$$p = \frac{RT}{V_m - b} - \frac{a}{V_m^2}$$
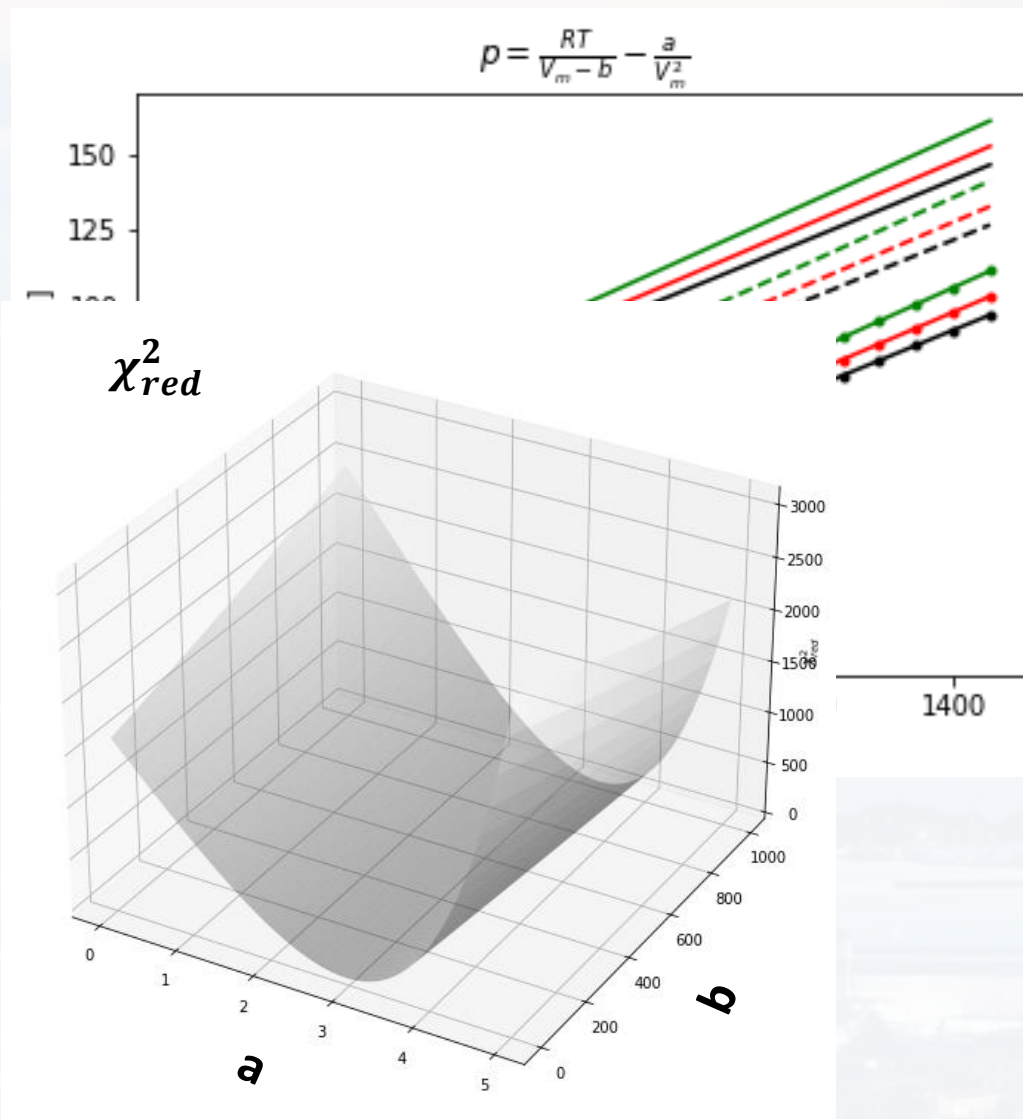
finding **a** and **b** of
a van-der-Waals gas

if critical points are not
accessible
→ fitting curve, finding **a** and **b**

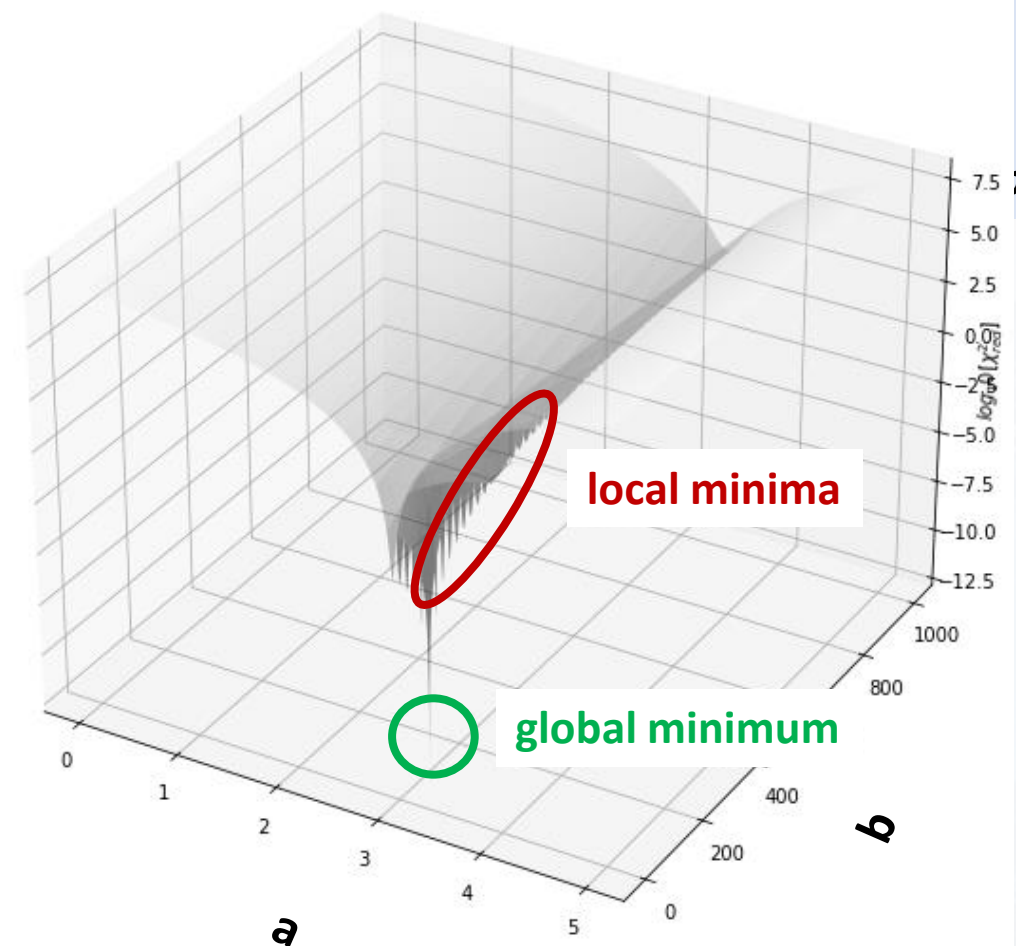Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)



$$p = \frac{RT}{V_m - b} - \frac{a}{V_m^2}$$

$\log(\chi_{red}^2)$

$\chi_{red}^2$

finding **a** and **b** of

local minima

global minimum

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

Often, the extreme of the objective function is subject to **constrains**

cross entropy $$S = -\sum_i p(true)_i \cdot \ln p(model)_i$$ constrain: $$\sum_i p_i = 1$$

→ Lagrangian Multipliers and variational calculus

→ mathematically: *Free Energy like term = Energy like term – Entropy term*

examples: - KL divergence
- Lasso method (linear regression)
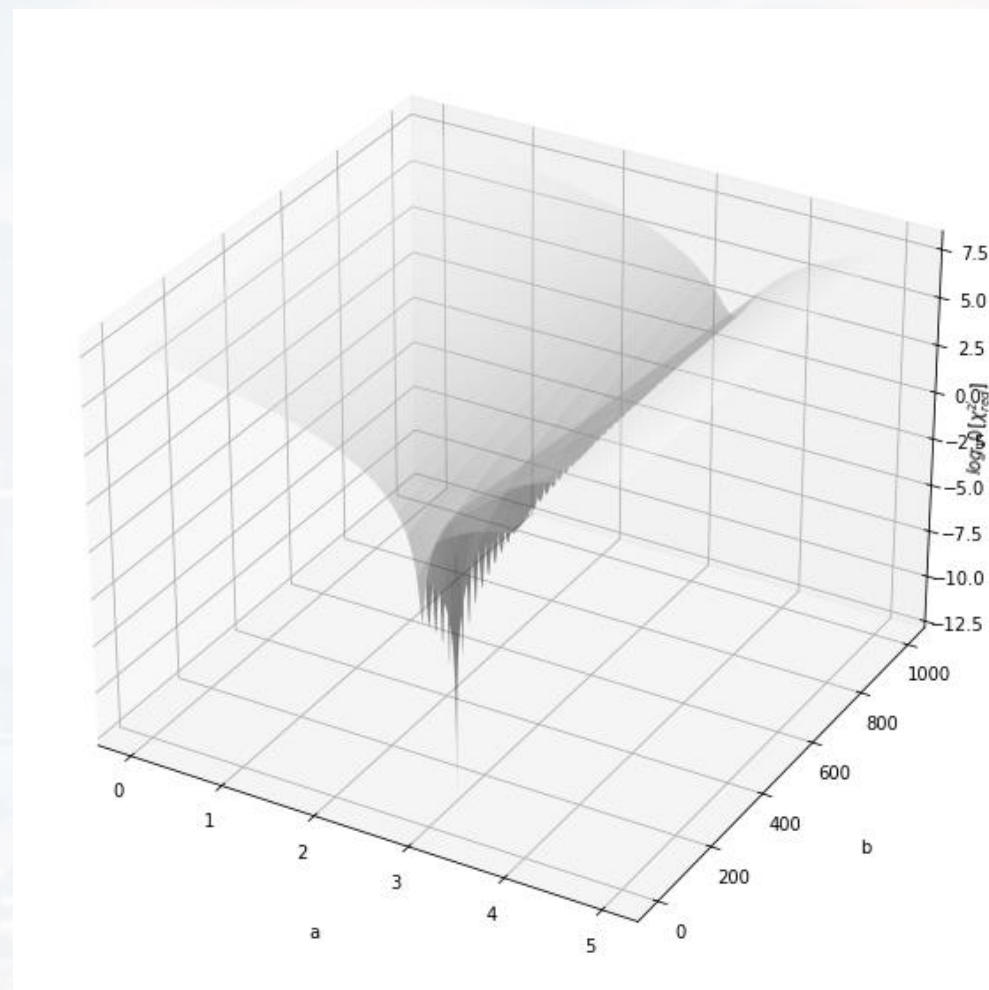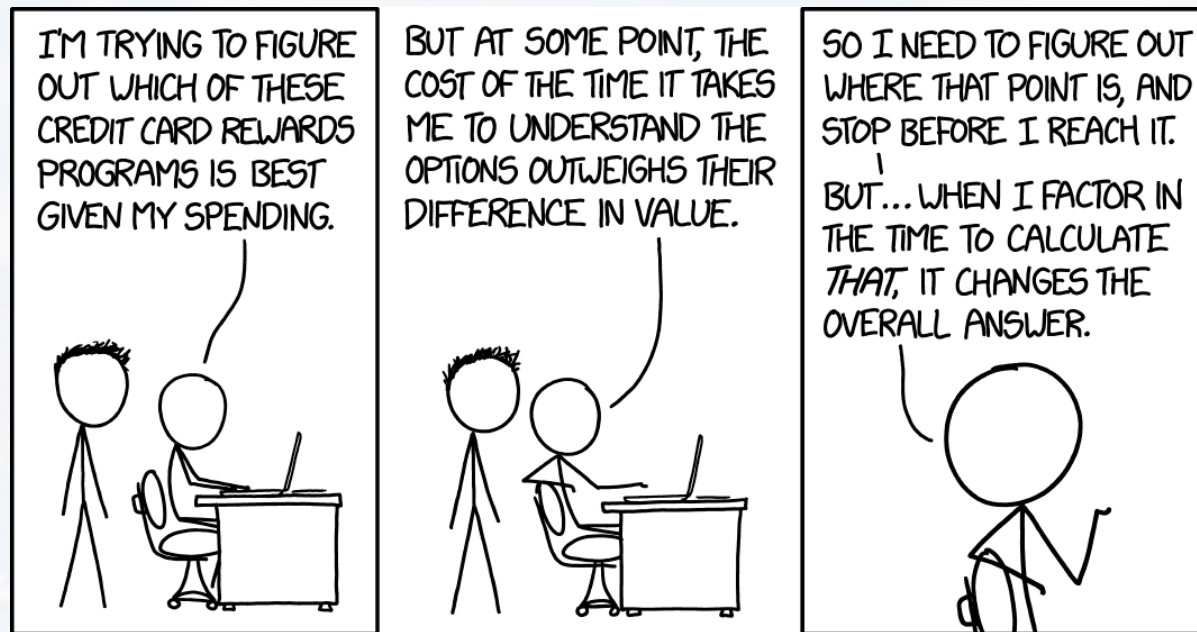- actual energy → Boltzmann distribution

etc

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

These functions are very complicated, not analytical at all

Outline

main application: **ANN!**

**Vanilla**

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

$y(x(t = 1))$

$\dfrac{dy}{dx}\bigg|_{x(t=1)}$

$x(t = 1)$    $x(t = 2)$

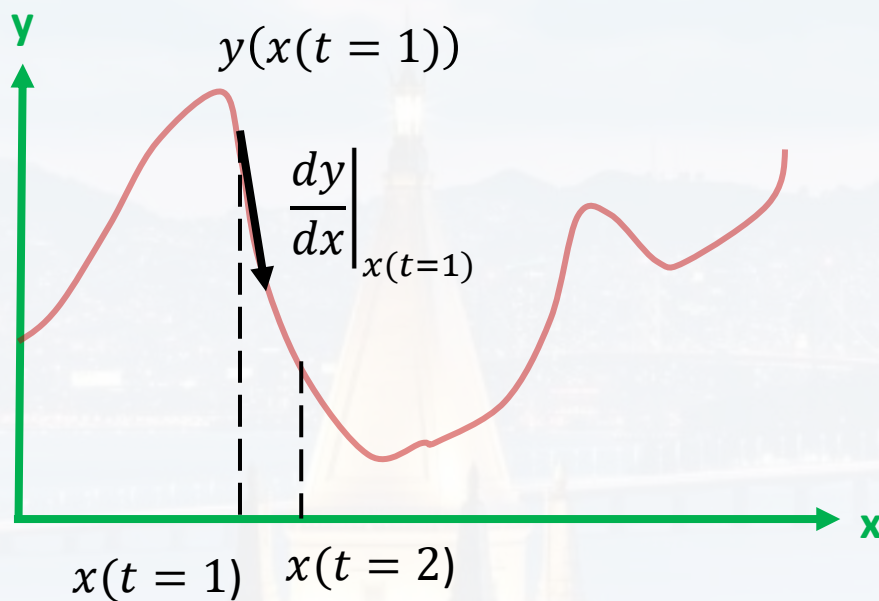$$x(t = 2) = x(t = 1) - \varepsilon \frac{dy}{dx}\bigg|_{x(t=1)}$$

$\varepsilon > 0$

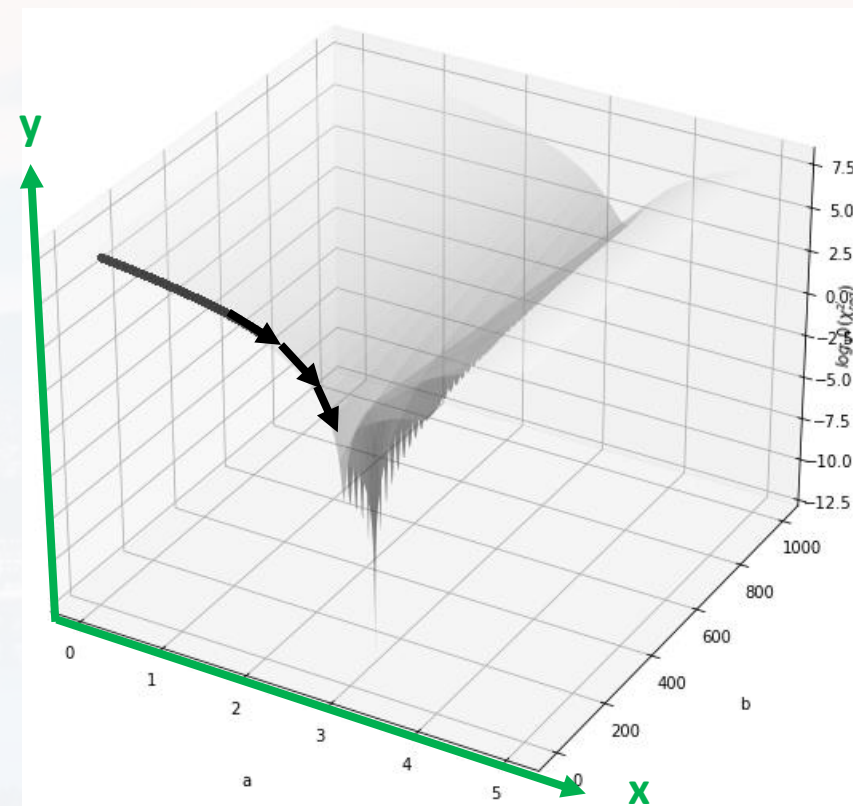**Vanilla**

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$



$$x(t=3) = x(t=2) - \varepsilon \frac{dy}{dx}\bigg|_{x(t=2)}$$

$\varepsilon > 0$

**Vanilla**

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$



$$\frac{dy}{dx}\bigg|_{x(t=3)}$$

$x(t=2)$    $x(t=3)$

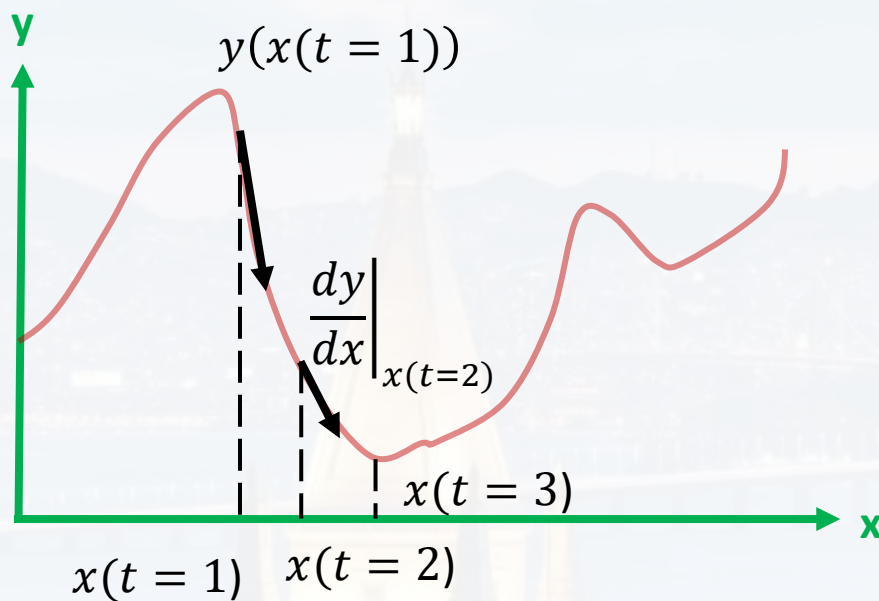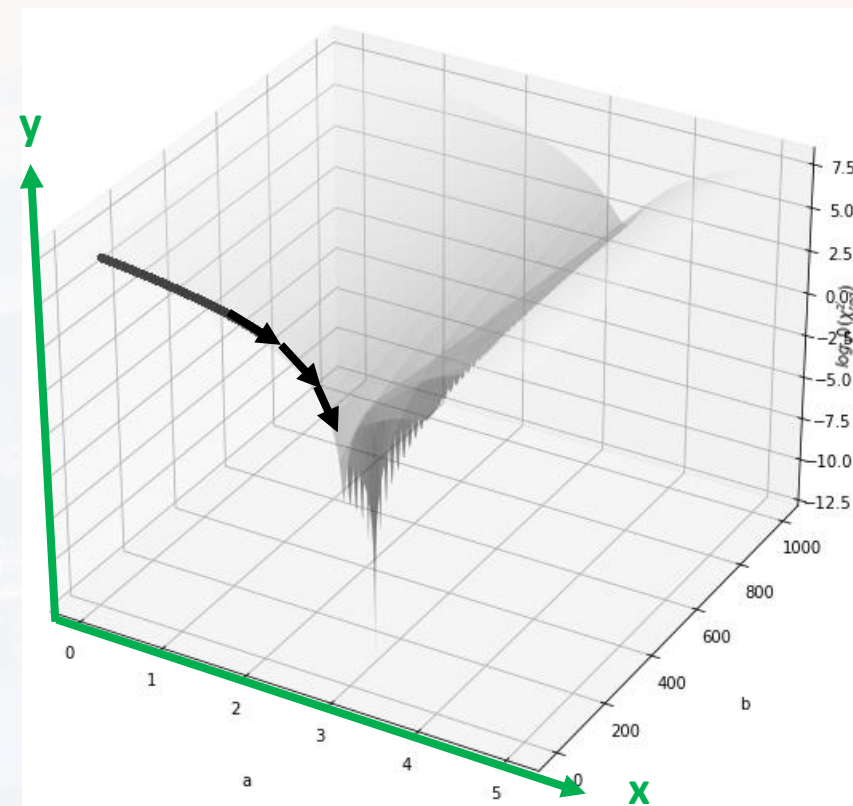$$x(t=4) = x(t=3) - \varepsilon \frac{dy}{dx}\bigg|_{x(t=3)}$$

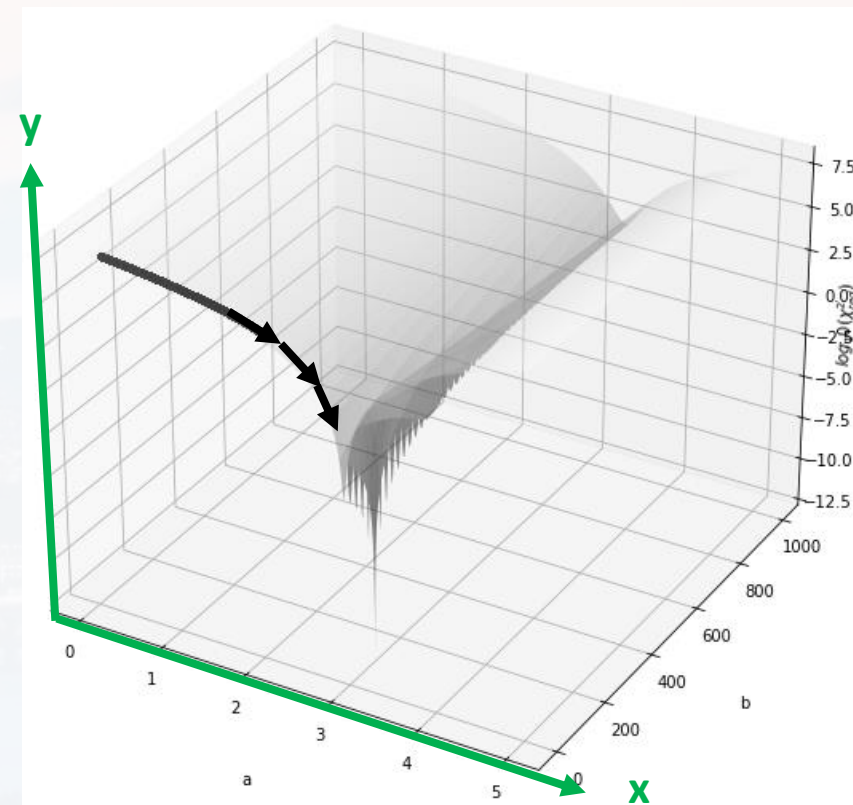$\varepsilon > 0$

**Vanilla**

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$



$x(t = 3)$

$$x(t = 4) = x(t = 3) - \varepsilon \frac{dy}{dx}\bigg|_{x(t=3)}$$

$\varepsilon > 0$

$$\frac{dy}{dx}\Bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$



$y(x(t=1))$

$\dfrac{dy}{dx}\Bigg|_{x(t=1)}$

$x(t=2)$  $x(t=1)$

$$x(t=2) = x(t=1) \; - \; \varepsilon \frac{dy}{dx}\Bigg|_{x(t=1)}$$

$\varepsilon > 0$

**Vanilla**

$$\frac{dy}{dx_1}\bigg|_{x_1(0)} \approx \frac{y(x_1(0) + \Delta x_1) - y(x_1(0) - \Delta x_1)}{2\Delta x_1}$$

$$\frac{dy}{dx_2}\bigg|_{x_2(0)} \approx \frac{y(x_2(0) + \Delta x_2) - y(x_2(0) - \Delta x_2)}{2\Delta x_2}$$

**Vanilla**

$$\left.\frac{dy}{dx_1}\right|_{x_1(0)} \approx \frac{y(x_1(0)+\Delta x_1)-y(x_1(0)-\Delta x_1)}{2\Delta x_1}$$

$$\left.\frac{dy}{dx_2}\right|_{x_2(0)} \approx \frac{y(x_2(0)+\Delta x_2)-y(x_2(0)-\Delta x_2)}{2\Delta x_2}$$

$$\vdots$$

$$\left.\frac{dy}{dx_i}\right|_{x_i(0)} \approx \frac{y(x_i(0)+\Delta x_i)-y(x_i(0)-\Delta x_i)}{2\Delta x_i}$$

$$\vdots$$

$$\left.\frac{dy}{dx_N}\right|_{x_N(0)} \approx \frac{y(x_N(0)+\Delta x_N)-y(x_N(0)-\Delta x_N)}{2\Delta x_N}$$

$$\begin{pmatrix} \left.\dfrac{dy}{dx_1}\right|_{x_1(0)} \\ \\ \dots \\ \\ \left.\dfrac{dy}{dx_i}\right|_{x_i(0)} \\ \\ \dots \\ \\ \left.\dfrac{dy}{dx_N}\right|_{x_N(0)} \end{pmatrix} = grad(y)_x$$

Outline

- **Gradient Descent**

    - Vanilla
    - Learning Rate Schedule
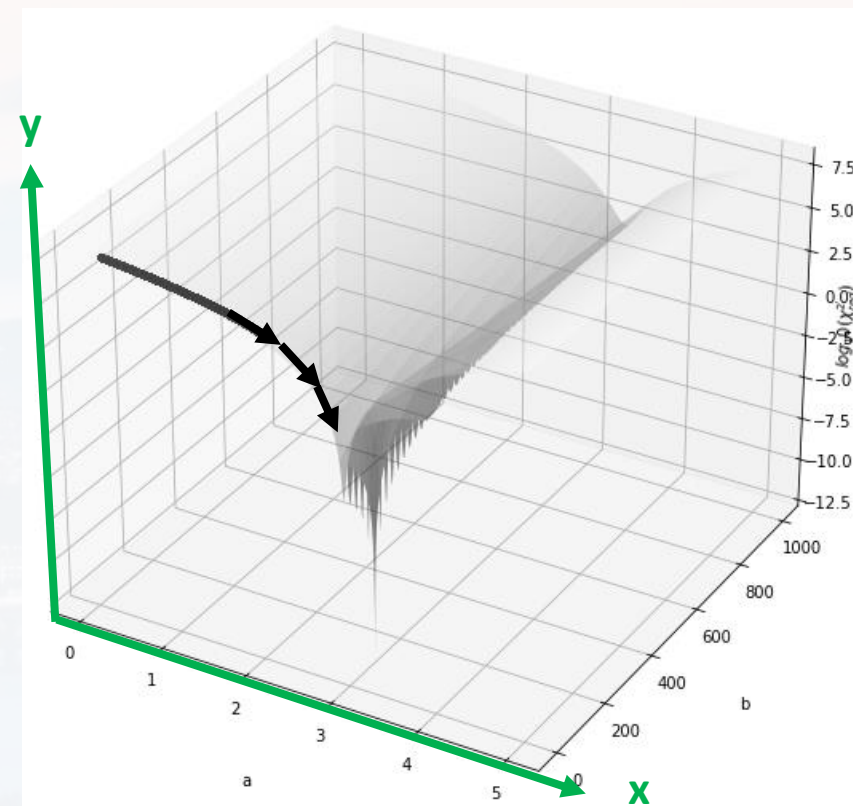    - Momentum
    - L1 and L2
    - More Finetuning

$$\left.\frac{dy}{dx}\right|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

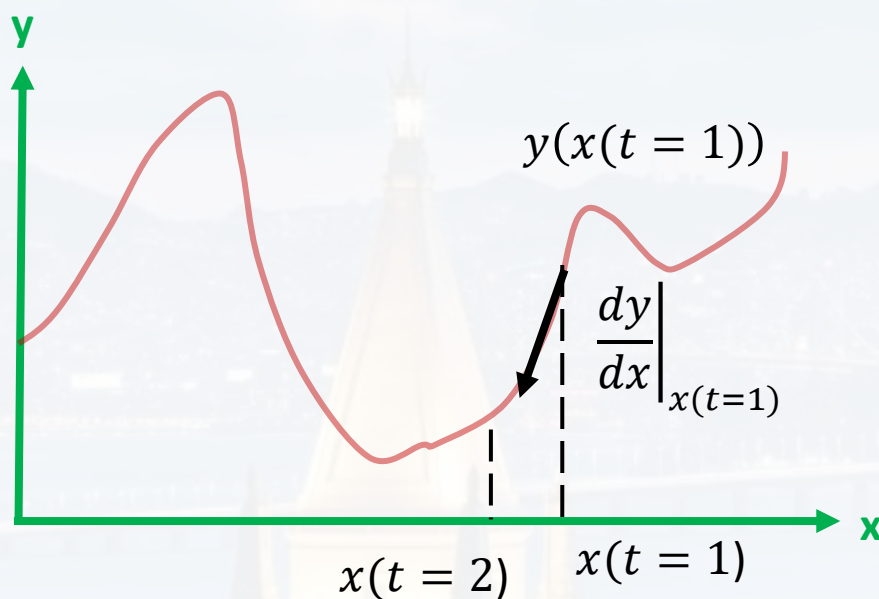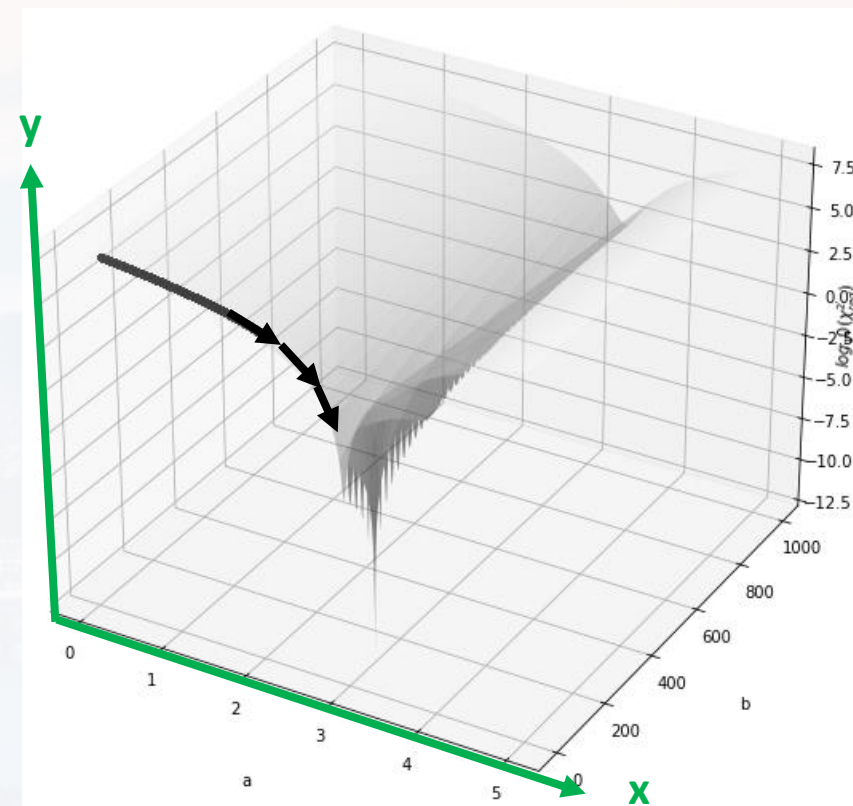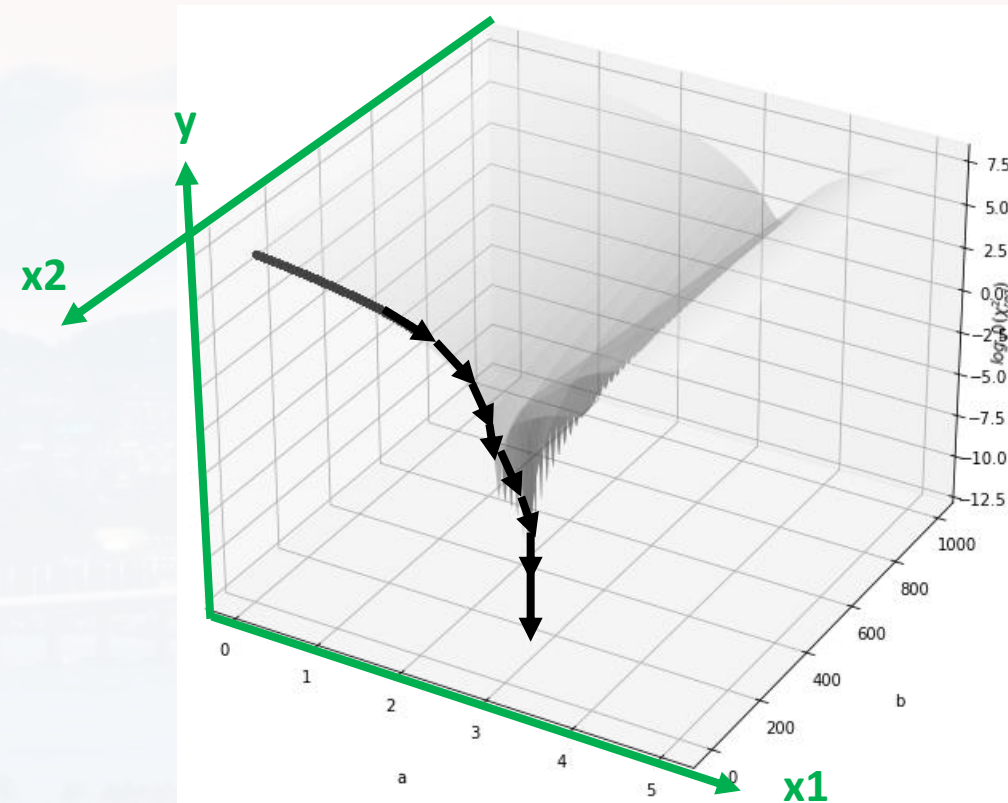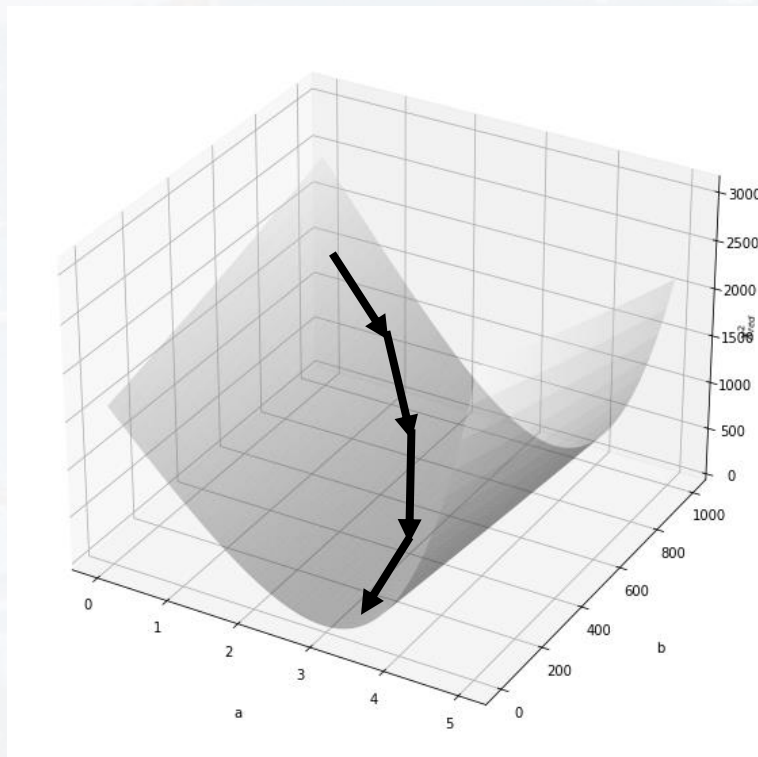$$x(t+1) = x(t) - \boldsymbol{\varepsilon}\left.\frac{dy}{dx}\right|_{x(t)}$$

**Learning Rate Schedule**

$$\boldsymbol{\varepsilon} > 0 \qquad \text{called } \textit{learning rate}$$

$$\Delta x = - \boldsymbol{\varepsilon}\left.\frac{dy}{dx}\right|_{x(t)} \qquad \text{defines how large the leap } \Delta x \text{ is}$$

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

$$x(t + 1) = x(t) - \varepsilon \frac{dy}{dx}\bigg|_{x(t)}$$

**Learning Rate Schedule**

$$\varepsilon > 0$$

called *learning rate*

$$\Delta x = - \varepsilon \frac{dy}{dx}\bigg|_{x(t)}$$

defines how large the leap $\Delta x$ is

$$\frac{dy}{dx}\Bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

$$x(t+1) = x(t) - \boldsymbol{\varepsilon}\frac{dy}{dx}\Bigg|_{x(t)}$$

**Learning Rate Schedule**

$\boldsymbol{\varepsilon} > 0$    called *learning rate*

$$\Delta x = -\boldsymbol{\varepsilon}\frac{dy}{dx}\Bigg|_{x(t)}$$

defines how large
the leap $\Delta x$ is

… and so on…

→ smaller $\boldsymbol{\varepsilon}$ ?

y

x

$x(t=4)$

$x(t=5)$

$$\frac{dy}{dx}\Bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

$$x(t + 1) = x(t) - \boldsymbol{\varepsilon}\frac{dy}{dx}\Bigg|_{x(t)}$$

$\boldsymbol{\varepsilon} > 0$    called *learning rate*



$$\Delta x = -\boldsymbol{\varepsilon}\frac{dy}{dx}\Bigg|_{x(t)}$$    defines how large the leap $\Delta x$ is

… and so on…

→ smaller $\boldsymbol{\varepsilon}$ ?    Takes too long!

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

$$x(t+1) = x(t) - \boldsymbol{\varepsilon}\frac{dy}{dx}\bigg|_{x(t)}$$

learning rate as function of t:

$$\boldsymbol{\varepsilon} > 0 \qquad \text{called } \textit{learning rate}$$

$$\boldsymbol{\varepsilon}(t) = \frac{\boldsymbol{\varepsilon}_0}{1 + \kappa\, t} \qquad \text{decay rate } \kappa$$

$$\Delta x = -\boldsymbol{\varepsilon}\frac{dy}{dx}\bigg|_{x(t)} \qquad$$

defines how large the leap $\Delta x$ is

$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$

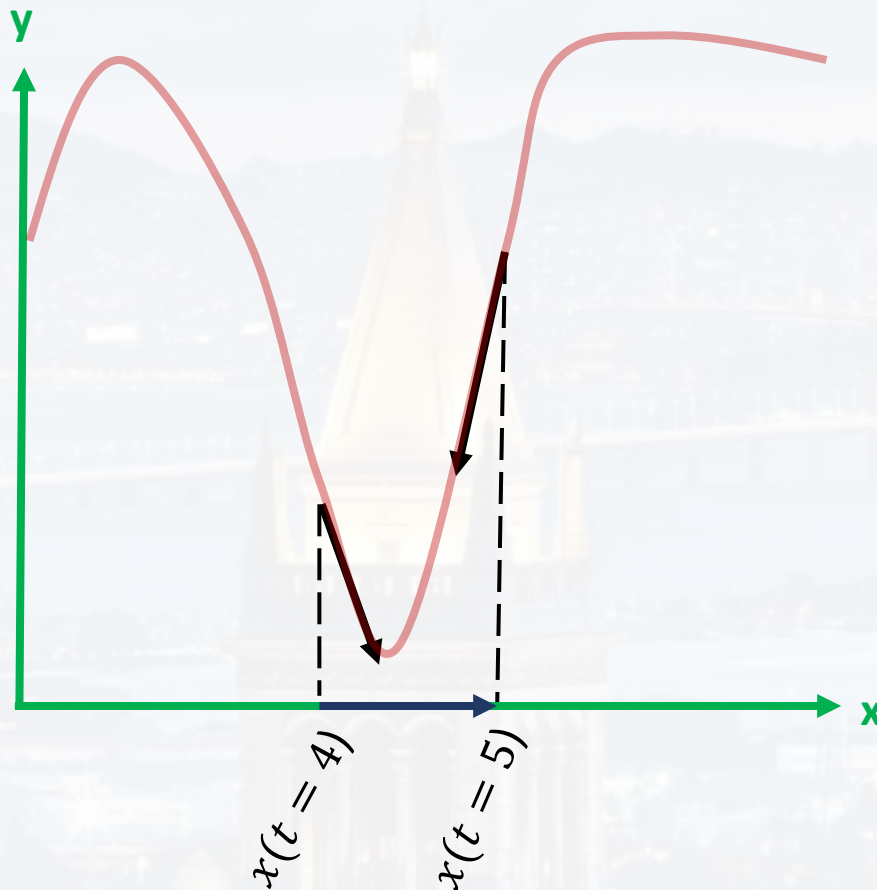$$x(t + 1) = x(t) - \varepsilon \frac{dy}{dx}\bigg|_{x(t)}$$

**Learning Rate Schedule**

learning rate as function of t:

$$\varepsilon > 0$$ called *learning rate*

$$\varepsilon(t) = \frac{\varepsilon_0}{1 + \kappa t} \quad \text{decay rate } \kappa$$

$$\Delta x = - \varepsilon \frac{dy}{dx}\bigg|_{x(t)}$$

defines how large the leap $\Delta x$ is

can also be a stepwise function (learning rate schedule)

learning rate as function of t:

$$\varepsilon(t) = \frac{\varepsilon_0}{1 + \kappa\, t}$$   decay rate $\kappa$

$$\Delta x = -\varepsilon \left.\frac{dy}{dx}\right|_{x(t)}$$   defines how large the leap $\Delta x$ is

can also be a stepwise function (learning rate schedule)



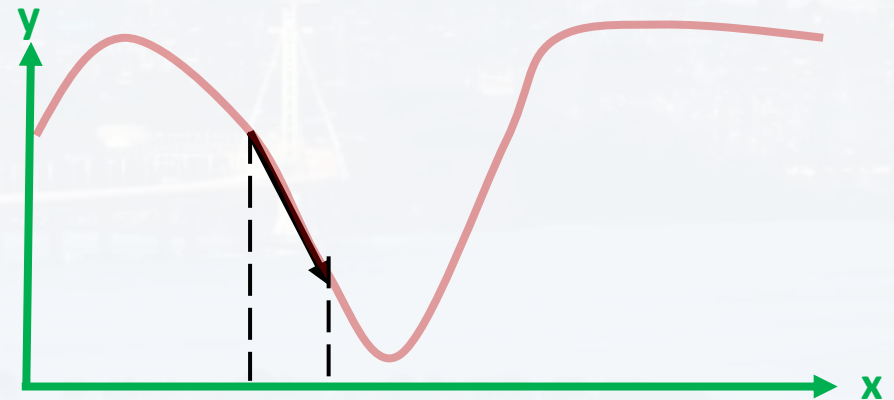$$\varepsilon \rightarrow \frac{\varepsilon}{\sqrt{grad(y)_x}}$$   adaptive gradient, aka **AdaGrad**

Outline

- The Problem

- **Gradient Descent**

  - Vanilla
  - Learning Rate Schedule
  - Momentum
  - L1 and L2
  - More Finetuning

**Momentum**



$$p = \frac{RT}{V_m - b} - \frac{a}{V_m^2}$$

even with AdaGrad and learning rate schedule
→ would get stuck in local minimum

need to roll over → **momentum**

**Momentum**

taking the **average** of *N* previous gradients
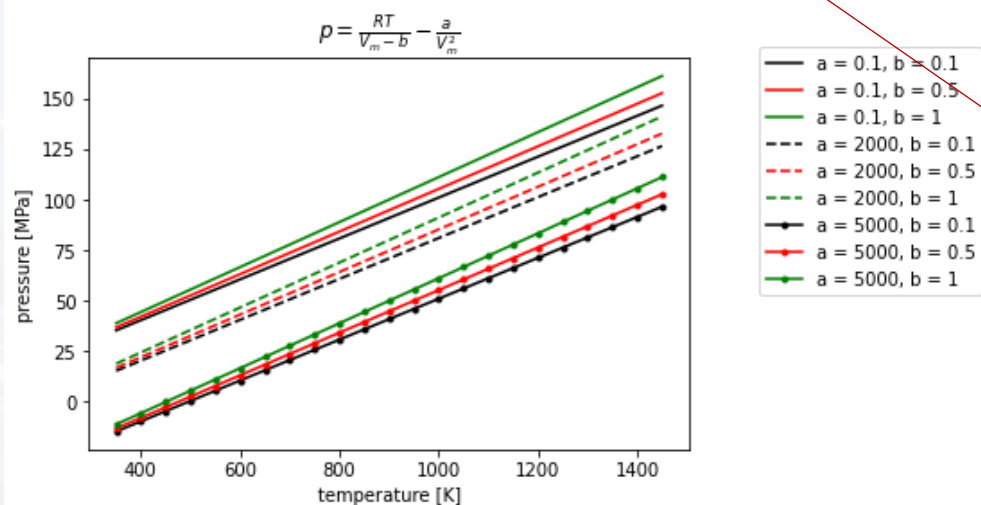
$$\langle grad(y)_{x(t)} \rangle = \frac{1}{N}[grad(y)_{x(t-1)} + grad(y)_{x(t-2)} +$$

$$\dots + grad(y)_{x(t-N)}]$$

but we want more recent gradients to contribute more than older gradients

→ **weighted average** with weighting factor $\boldsymbol{\mu_k}$

$$\langle grad(y)_{x(t)} \rangle = \sum_{k=t-N}^{t-1} \mu_k \cdot grad(y)_{x(k)}$$

Finding a clever way to adjust $\boldsymbol{\mu_k}$ during every iteration *t*

**weighted average** with weighting factor $\boldsymbol{\mu_k}$

Finding a clever way to adjust $\boldsymbol{\mu_k}$ during every iteration $t$

$$\langle grad(y)_{x(0)} \rangle = grad(y)_{x(0)} \qquad \mu_0 = (0,1)$$

$$\langle \boldsymbol{grad(y)_{x(1)}} \rangle = grad(y)_{x(1)} + \mu_0 \cdot grad(y)_{x(0)}$$

**weighted average** with weighting factor $\boldsymbol{\mu_k}$

Finding a clever way to adjust $\boldsymbol{\mu_k}$ during every iteration $t$

$$\langle grad(y)_{x(0)} \rangle = grad(y)_{x(0)} \qquad \mu_0 = (0,1)$$

$$\langle \boldsymbol{grad}(\boldsymbol{y})_{x(1)} \rangle = grad(y)_{x(1)} + \mu_0 \cdot grad(y)_{x(0)}$$

$$\langle \boldsymbol{grad}(\boldsymbol{y})_{x(2)} \rangle = grad(y)_{x(2)} + \boxed{\mu_0} [grad(y)_{x(1)}$$

$$+ \boxed{\mu_0} grad(y)_{x(0)}]$$

$$\boldsymbol{\mu_{k=2}} = \mu_0 \, \mu_0 = \mu_0^2$$

$$\langle \boldsymbol{grad}(\boldsymbol{y})_{x(3)} \rangle = grad(y)_{x(3)} + \boxed{\mu_0} [grad(y)_{x(2)} + \boxed{\mu_0} [grad(y)_{x(1)} + \boxed{\mu_0} \cdot grad(y)_{x(0)}]]$$
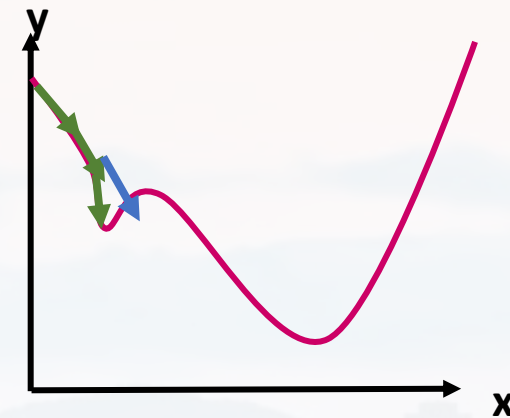
… and so on…

**weighted average** with weighting factor $\boldsymbol{\mu_k}$

$\mu_0 = (0,1)$    called ''momentum'

$$\langle \boldsymbol{grad}(y)_{x(3)} \rangle = grad(y)_{x(3)} +$$

$$\mu_0 \left[ grad(y)_{x(2)} + \mu_0 \left[ grad(y)_{x(1)} + \mu_0 \cdot grad(y)_{x(0)} \right] \right] \quad \text{... and so on...}$$

```python
class Optimizer:

    def __init__(self, learning_rate = 0.1, decay = 0, momentum = 0):
        self.learning_rate        = learning_rate
        self.decay                = decay
        self.current_learning_rate = learning_rate
        self.iterations           = 0
        self.momentum             = momentum
```

Outline

- The Problem

- **Gradient Descent**

  - Vanilla
  - Learning Rate Schedule
  - Momentum
  - L1 and L2
  - More Finetuning

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

Often, the extreme of the objective function is subject to **constrains**

sometimes we have some **prior knowledge** about the **independent variables**

recall: linear regression

finding best $\beta$ by

$$\min_{\beta}\left\{\frac{1}{N}\|Y - X\beta\|^2\right\}$$

now:
constrain: **encourages sparsity of $\beta$**

$$\min_{\beta}\left\{\frac{1}{N}\|Y - X\beta\|^2 + \lambda\|\beta\|^1\right\}$$

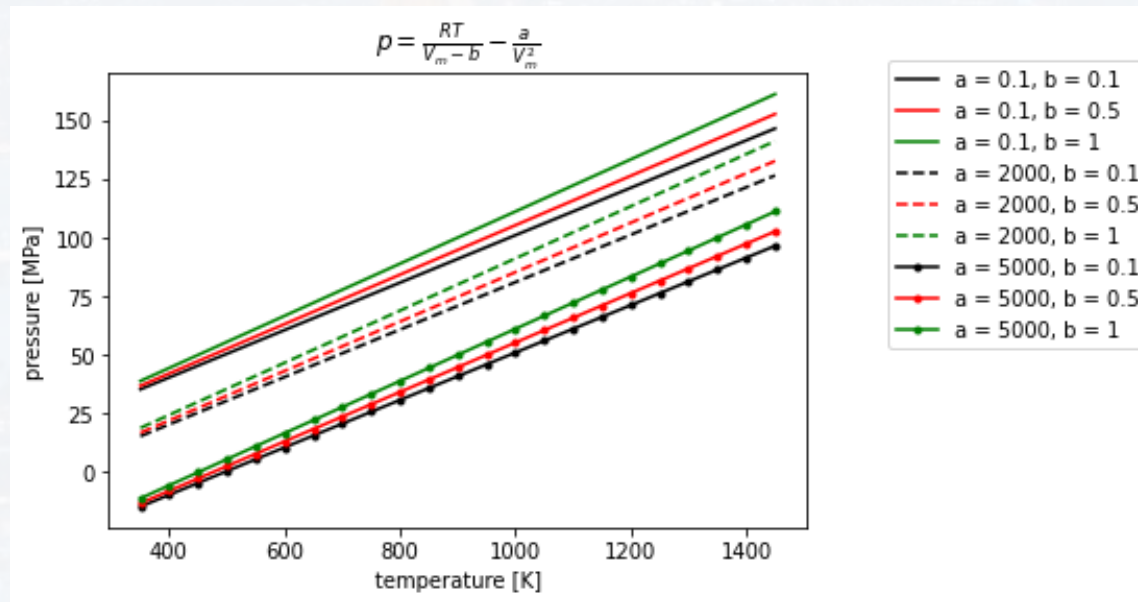$\lambda$    *Lagrangian Multiplier*

called **L1 regularization**, or LASSO

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

Often, the extreme of the objective function is subject to **constrains**

sometimes we have some **prior knowledge** about the **independent variables**

**L1 regularization**

L1 and L2



$$p = \frac{RT}{V_m - b} - \frac{a}{V_m^2}$$

We often have even hard constrains based on the laws of physics!

Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

Often, the extreme of the objective function is subject to **constrains**

sometimes we have some **prior knowledge** about the **independent variables**

recall: linear regression

finding best $\beta$ by

$$\min_{\beta} \left\{ \frac{1}{N} \|Y - X\beta\|^2 \right\}$$

now:

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y \longrightarrow \min_{\beta} \left\{ \frac{1}{N} \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \right\}$$
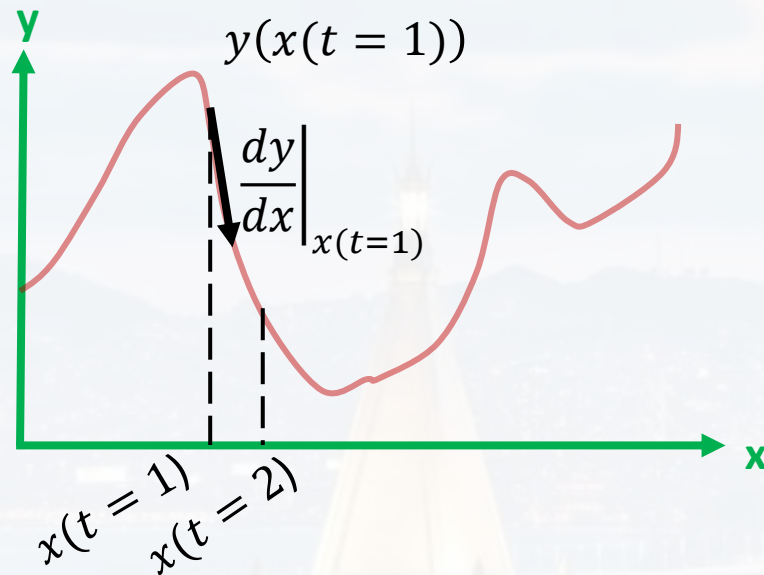
$\lambda$    *Lagrangian Multiplier*

called **L2 regularization**, or RIDGE  penalizes large $\beta$
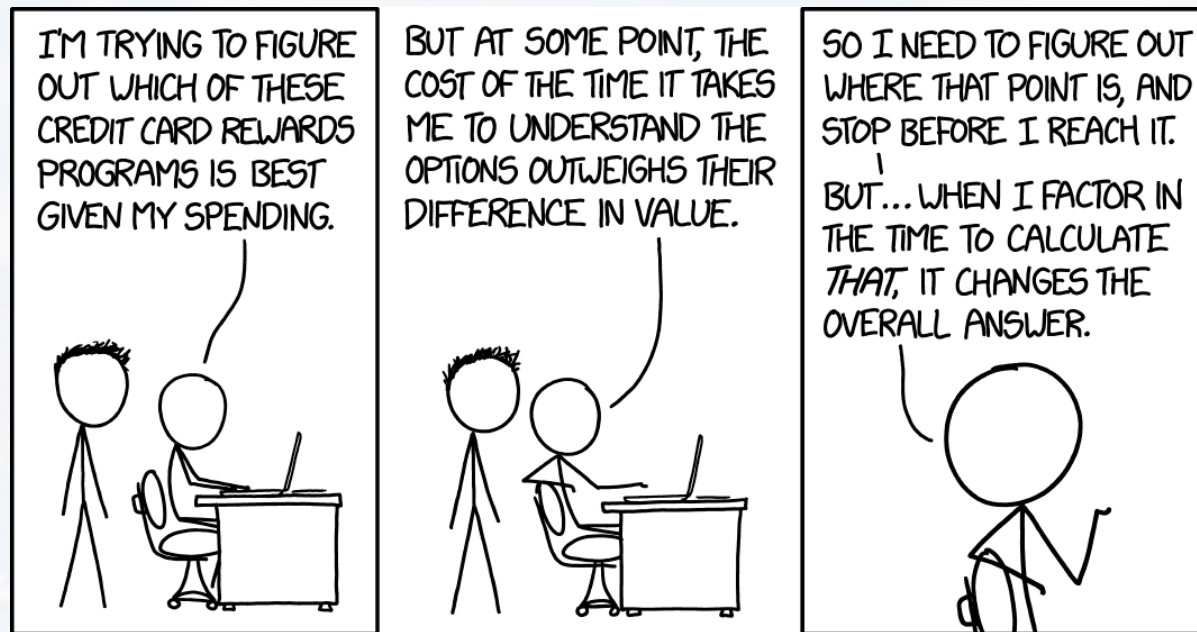
**L1 and L2 regularization**

L1 and L2

$$x(t = 2) = x(t = 1) - \varepsilon \frac{d[y + \lambda_1 \|x\|^1 + \lambda_2 \|x\|^2]}{dx}\bigg|_{x(t=1)}$$

$y(x(t = 1))$

$\dfrac{dy}{dx}\bigg|_{x(t=1)}$

$x(t = 1)$

$x(t = 2)$

- gradient descent does not stop if values for x are too large and prefers sparsity

- note: the derivative of $\|x\|^1$ returns the sign (i. e. direction)

- usually $\lambda \ll \|x\|^n$

- will be important for ANNs later

## Outline

- The Problem

- **Gradient Descent**

    - Vanilla
    - Learning Rate Schedule
    - Momentum
    - L1 and L2

- More Finetuning

Vanilla Gradient Descent → **S**tochastic **G**radient **D**escent

Learning Rate Schedule, L1, L2

Momentum

**different scaling for all different directions**

$$\varepsilon \ \rightarrow \ \frac{\varepsilon}{\sqrt{grad(y)_x}}$$

adaptive gradient, aka **AdaGrad**

Adding a decay factor to the sum of gradient squared (similar to momentum),
aka **R**oot **M**ean **S**quare **Prop**agation **RMSProp**

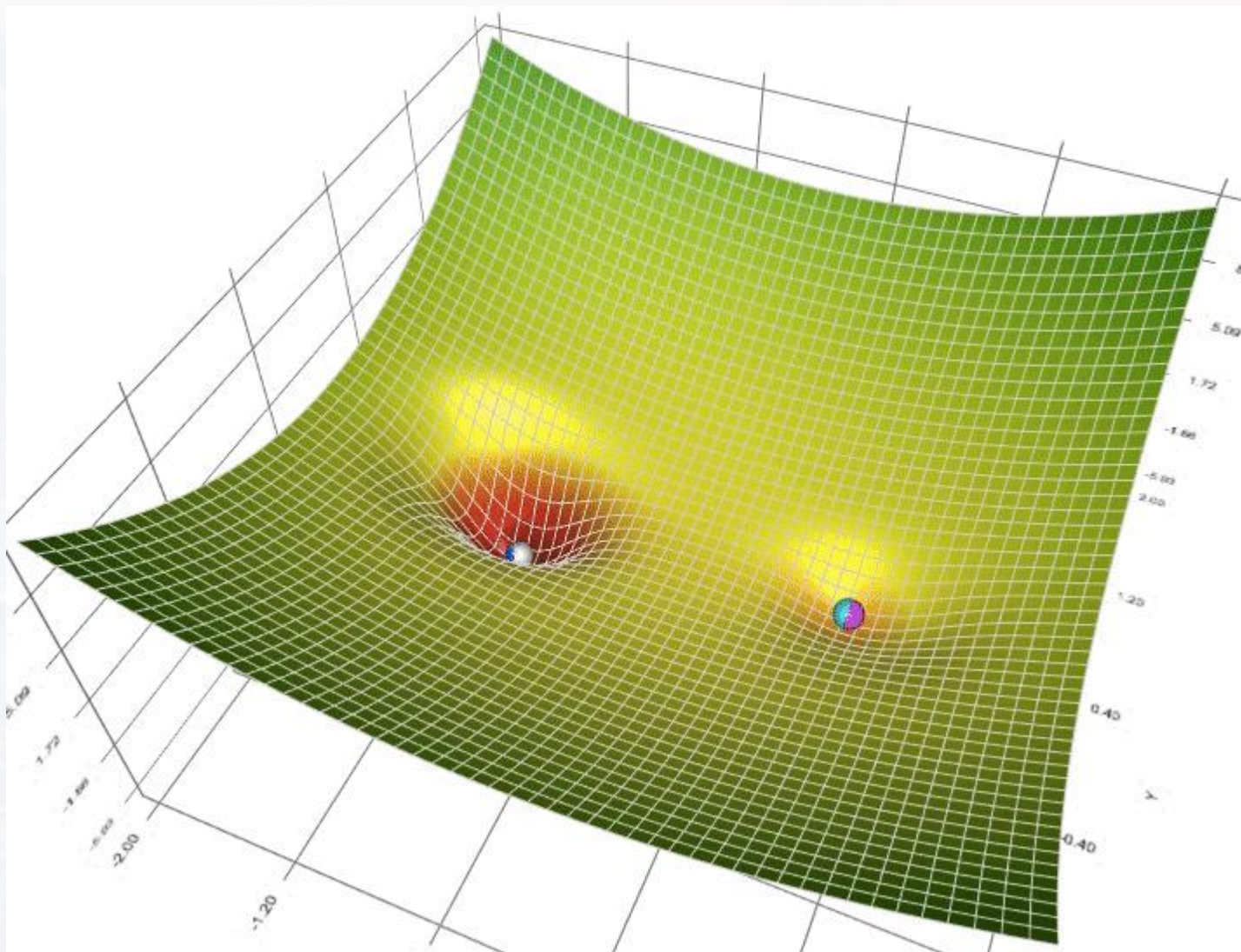all combined:
**Ada**ptive **M**oment Estimation
aka **Adam**

TowardsDataScience

**Lili Jiang**

**More Fine Tuning**

**gradient descent (cyan), momentum (magenta), AdaGrad (white), RMSProp (green), Adam (blue)**
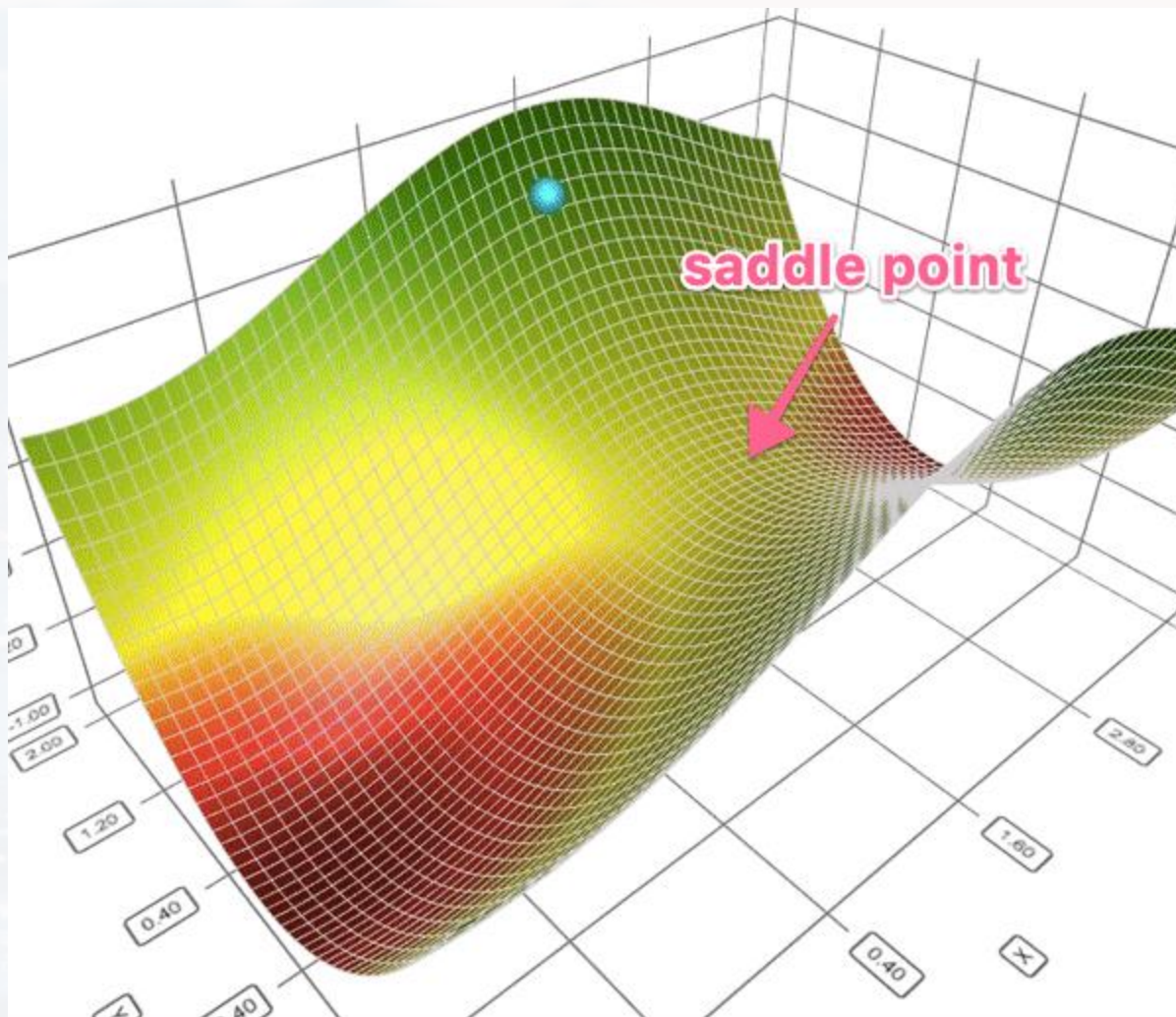
TowardsDataScience

**Lili Jiang**

**More Fine Tuning**

**gradient descent (cyan),
momentum (magenta),
AdaGrad (white),
RMSProp (green),
Adam (blue)**

**Thank you very much for your attention!**