**Lecture 06:**

**Optimization**

Markus Hohle

University California, Berkeley

**Machine Learning Algorithms**

MSSE 277B, 3 Units

# Course Map

**classic ML tools & algorithms**

**ANNs/AI/Deep Learning**

Outline

- **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:

$$\chi^2_{red} = \frac{1}{N-p-1} \sum_{i=1}^{N} \frac{(\hat{y}(model)_i - y_i)^2}{\sigma_i^2}$$

**classification**



maximize: accuracy

minimize: cross entropy

$$S = -\sum_i p(true)_i \cdot \ln p(model)_i$$

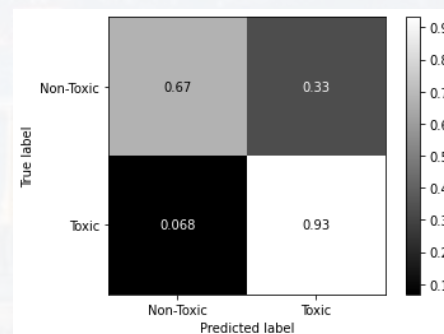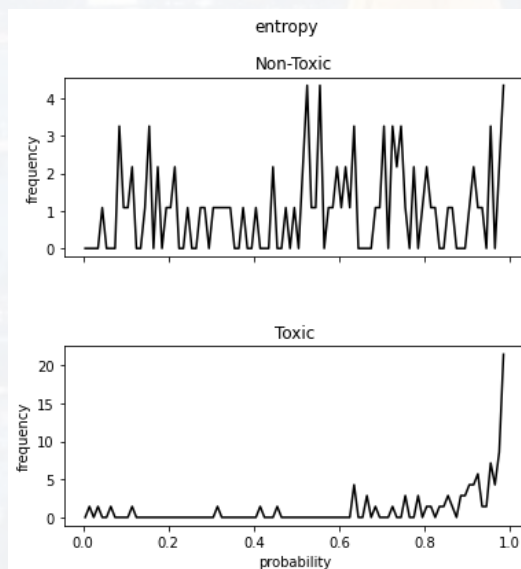Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

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

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

→ Lagrangian Multipliers and variational calculus

→ mathematically: $\qquad$ *Free Energy like term = Energy like term – Entropy term*

$\qquad\qquad$ examples: $\qquad\qquad$ - **E**vidence **L**ower **Bo**und
$\qquad\qquad\qquad\qquad\qquad\qquad$ - Lasso method (linear regression)
$\qquad\qquad\qquad\qquad\qquad\qquad$ - actual energy → Boltzmann distribution
$\qquad\qquad$ etc

# **Optimization**:

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

These functions are very complicated, not analytical ( = no mathematical equation) at all

<u>two most common approaches:</u>
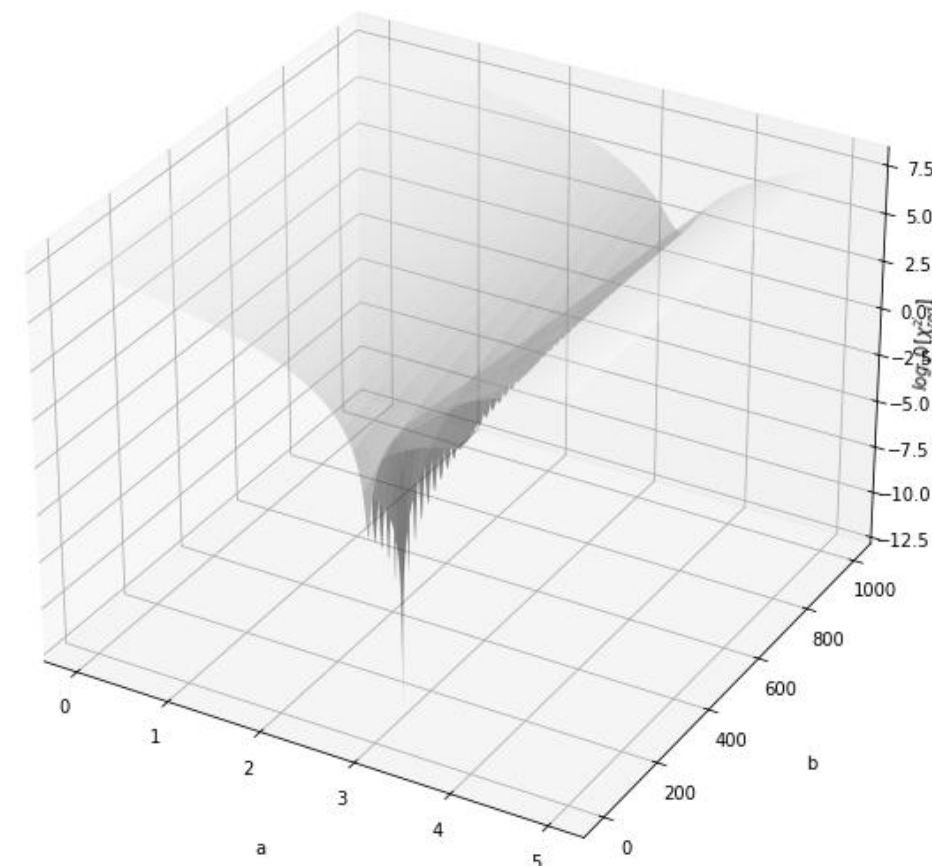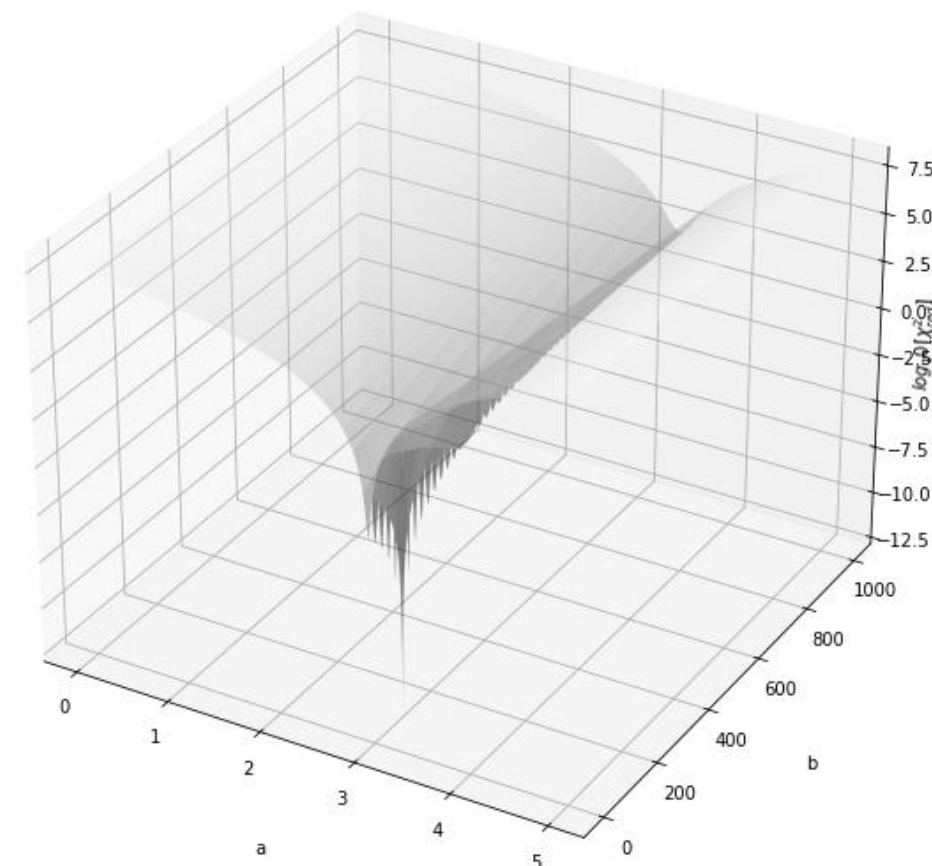
- gradient descent
- simulated annealing

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

These functions are very complicated, not analytical ( = no mathematical equation) at all

two most common approaches:

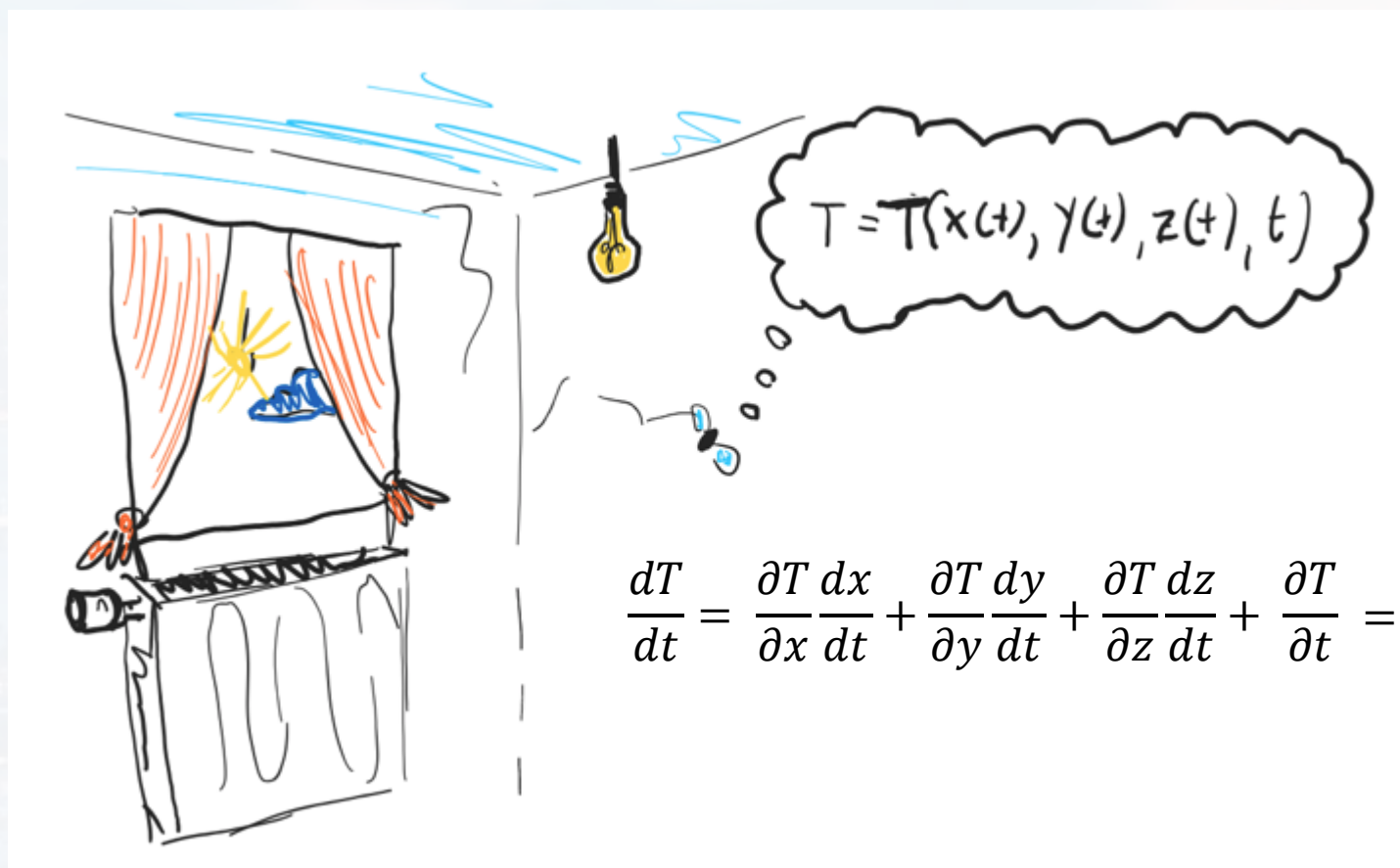**- gradient descent**
- simulated annealing

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

**gradient descent**

temperature profile in space and time

*T*: temperature

$$T = T(x(t), y(t), z(t), t)$$

$$\frac{dT}{dt} = \frac{\partial T}{\partial x}\frac{dx}{dt} + \frac{\partial T}{\partial y}\frac{dy}{dt} + \frac{\partial T}{\partial z}\frac{dz}{dt} + \frac{\partial T}{\partial t} = \begin{pmatrix} \frac{\partial T}{\partial x} \\ \frac{\partial T}{\partial y} \\ \frac{\partial T}{\partial z} \end{pmatrix} \circ \begin{pmatrix} \frac{dx}{dt} \\ \frac{dy}{dt} \\ \frac{dz}{dt} \end{pmatrix} + \frac{\partial T}{\partial t}$$

$$= grad(T) \circ \vec{v}$$

If *T* doesn't change with time!

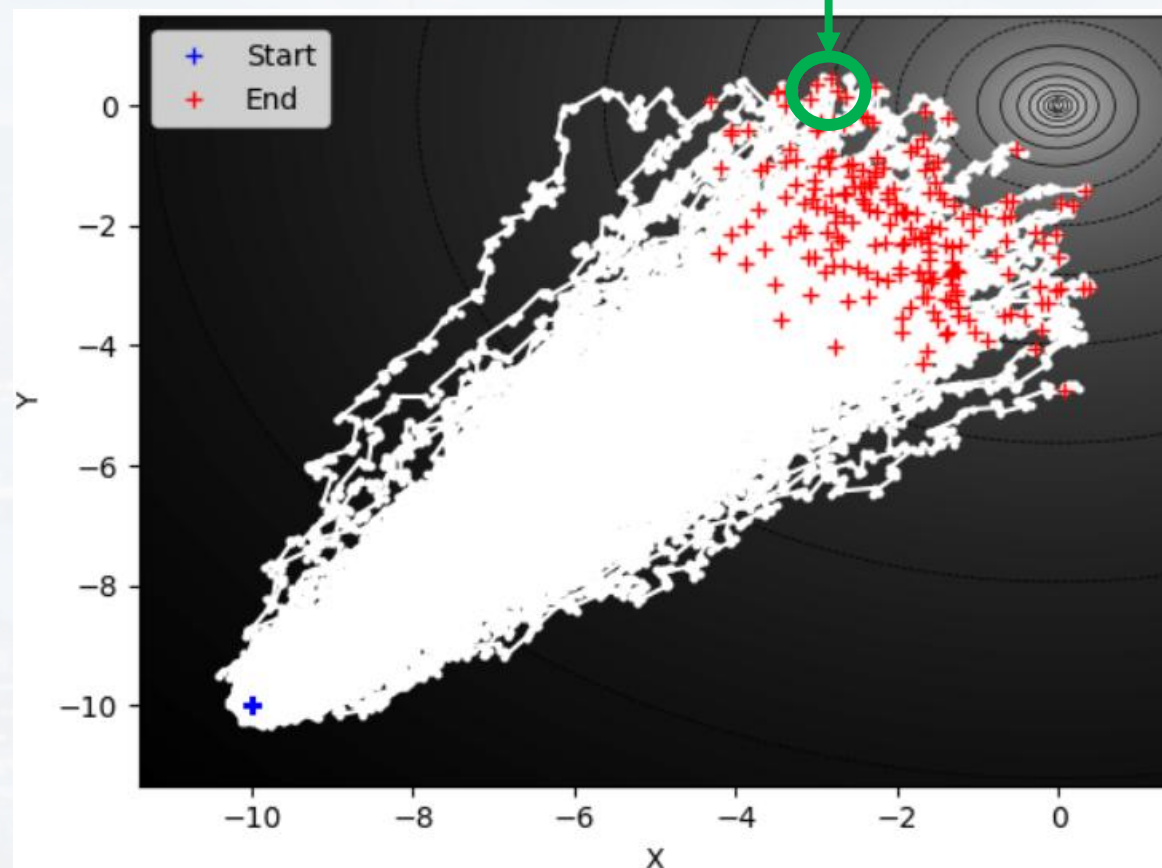Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)
→ extreme of an objective function

**gradient descent**

concentration profile in space and time

E. Coli

$c$: concentration



$$\frac{dc}{dt} = \frac{\partial c}{\partial x}\frac{dx}{dt} + \frac{\partial c}{\partial y}\frac{dy}{dt} + \frac{\partial c}{\partial z}\frac{dz}{dt} + \frac{\partial c}{\partial t}$$
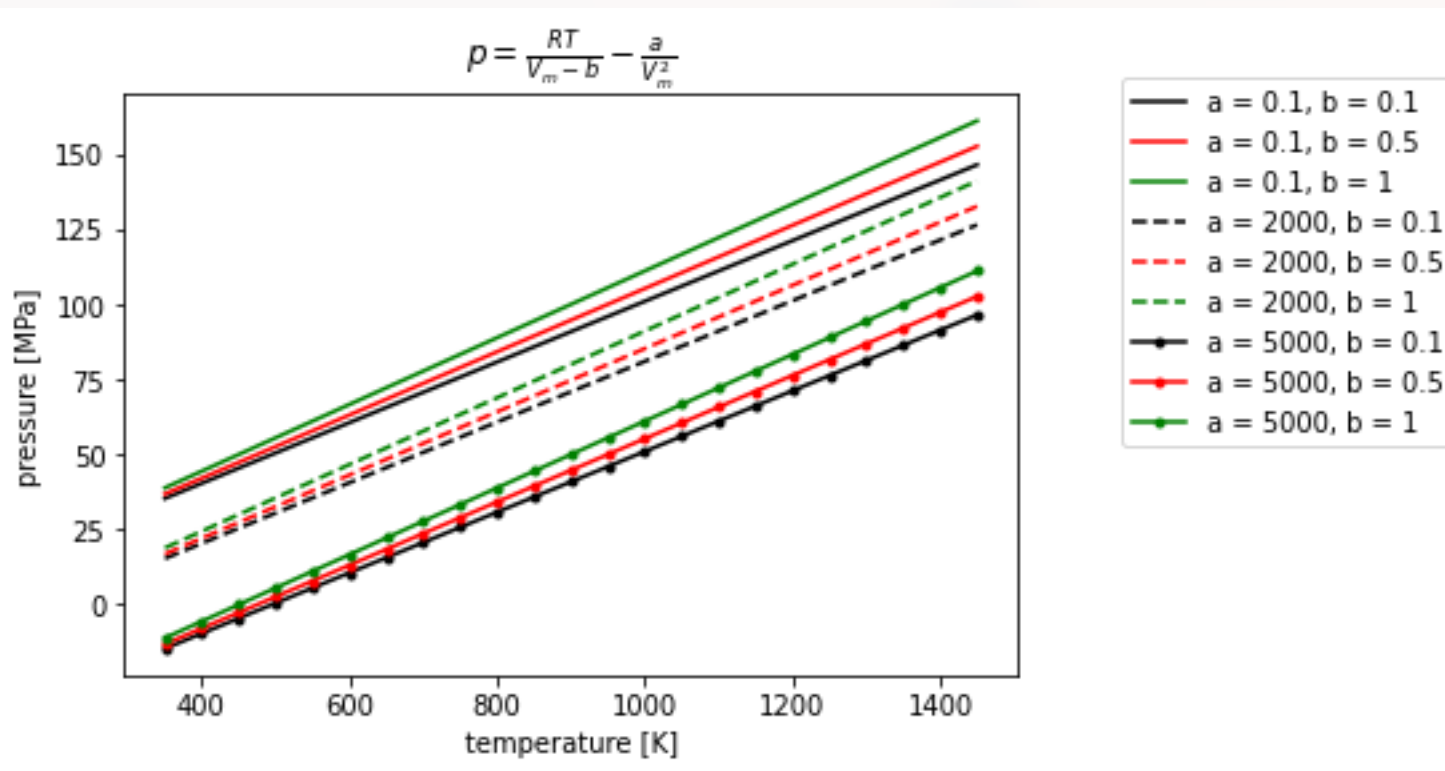
$$= \begin{pmatrix} \dfrac{\partial c}{\partial x} \\ \dfrac{\partial c}{\partial y} \\ \dfrac{\partial c}{\partial z} \end{pmatrix} \circ \begin{pmatrix} \dfrac{dx}{dt} \\ \dfrac{dy}{dt} \\ \dfrac{dz}{dt} \end{pmatrix} + \frac{\partial c}{\partial t}$$

$$= grad(c) \circ \vec{v}$$

If $c$ doesn't change with time!

# Optimization:

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

**gradient descent**



finding ***a*** and ***b*** of a van-der-Waals gas

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

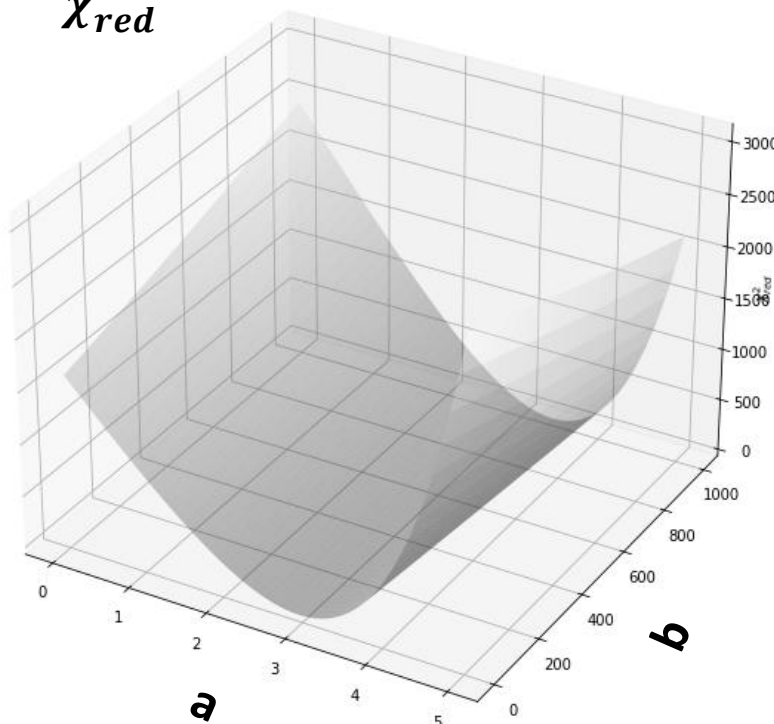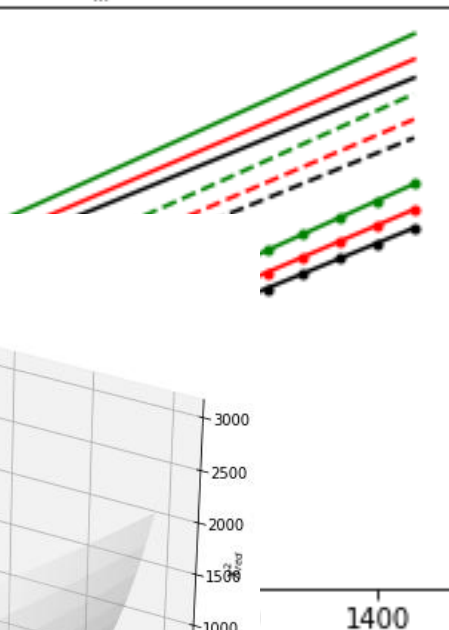Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

**gradient descent**

$$\chi^2_{red} = \frac{1}{N - p - 1} \sum_{i=1}^{N} \frac{(\hat{y}(model)_i - y_i)^2}{\sigma_i^2}$$

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

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

$\log(\chi^2_{red})$

$\chi^2_{red}$



local minima

global minimum

$a_{best}$   $b_{best}$

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

These functions are very complicated, not analytical ( = no mathematical equation) at all

two most common approaches:

    - gradient descent
    **- simulated annealing**

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

**simulated annealing**

"energy E(x,y)"

x

$\Delta E(x,y)$

y

If $\Delta E(x,y)$ is **negative**:
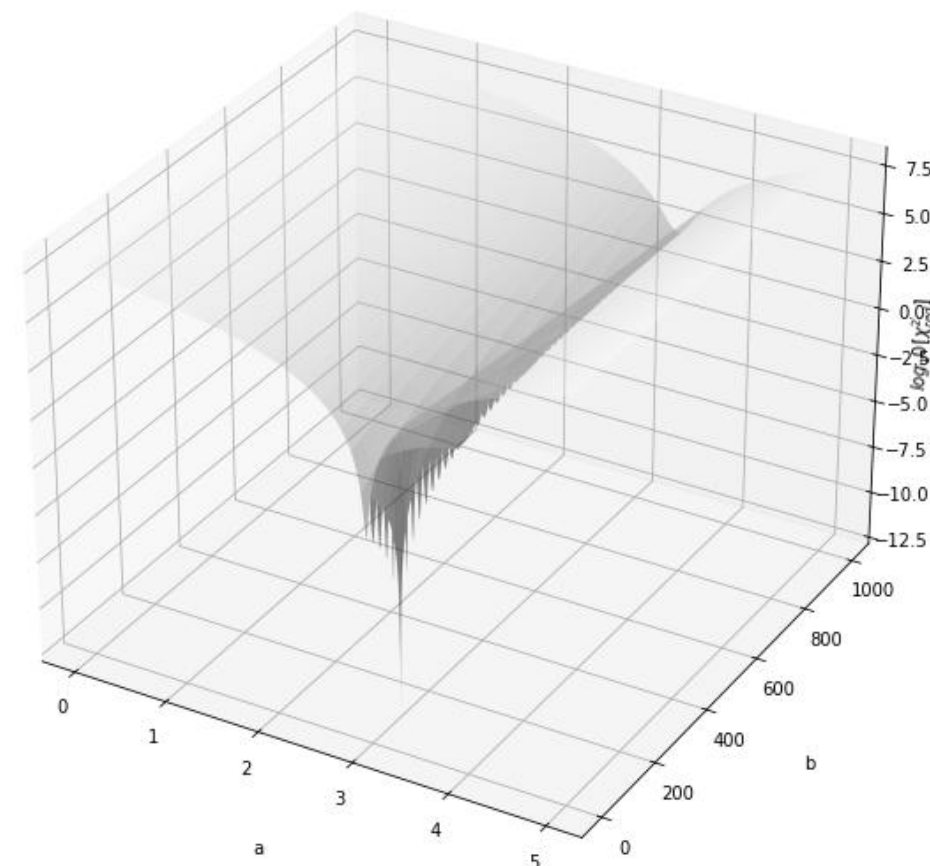→ **always move**
(a ball always rolls down the hill)

If $\Delta E(x,y)$ is **positive**:
→ calculate the **probability to move**
→ leaves some chance to escape local minimum

$T$: temperature

Boltzmann factor

$$p_{move} \sim \exp\left[-\frac{\Delta E(x,y)}{T}\right]$$

# Optimization:

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

**simulated annealing**



If $\Delta E(x, y)$ is **negative**:
→ **always move**
(a ball always rolls down the hill)

If $\Delta E(x, y)$ is **positive**:
→ calculate the **probability to move**
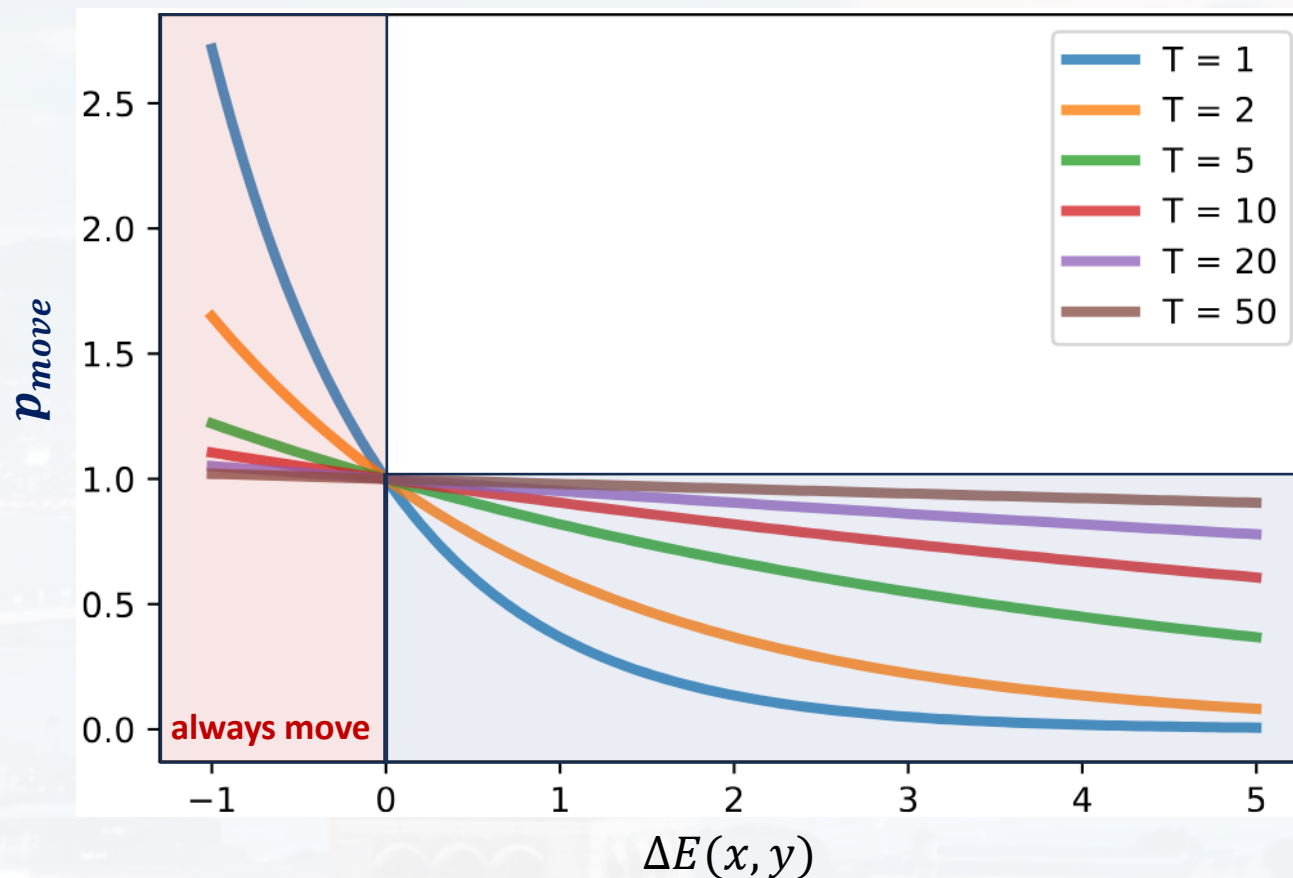→ leaves some chance to escape local minimum

$T$: temperature

Boltzmann factor

$$p_{move} \sim \exp\left[-\frac{\Delta E(x, y)}{T}\right]$$

**slowly reducing T → making larger jumps ($\Delta E(x, y)$) less likely over time**

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

**simulated annealing**

**Metropolis (Chem 273):**

1) suggest a random move $\Delta x$ and $\Delta y$

2) calculate $\Delta E(x, y)$ based on $\Delta x$ and $\Delta y$

3) move or not:

    a) move if $\Delta E(x, y) < 0$

    b) if $\Delta E(x, y) > 0$
- draw a **random number $\rho$** from a **uniform distribution** in the interval $(\mathbf{0}, \mathbf{1})$
- move if $\rho < \exp\left[-\frac{\Delta E(x,y)}{T}\right]$

4) reduce $T$ and repeat

Outline

- The Problem

- **Gradient Descent**

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

main application: **ANN!**
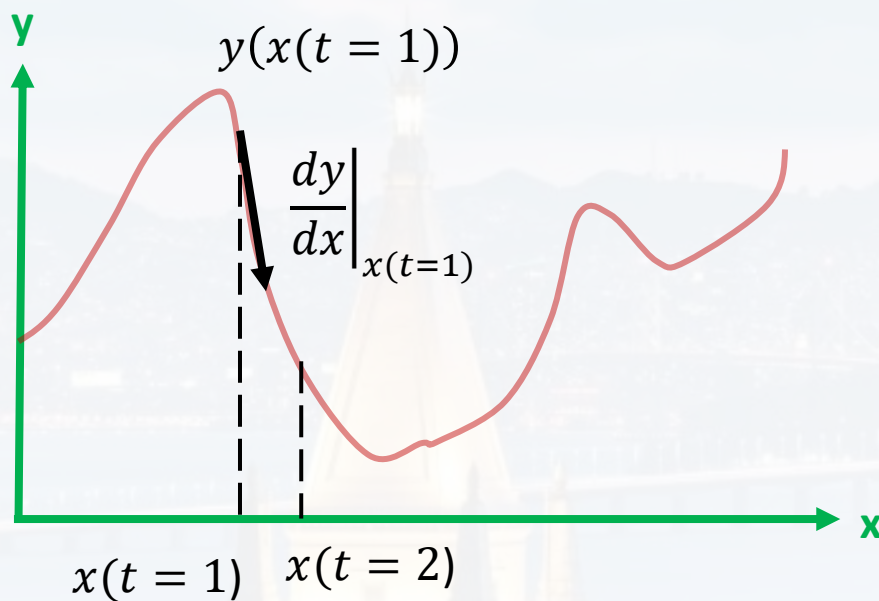
$$\frac{dy}{dx}\bigg|_{x_0} \approx \frac{y(x_0 + \Delta x) - y(x_0 - \Delta x)}{2\Delta x}$$



$y(x(t=1))$

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

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

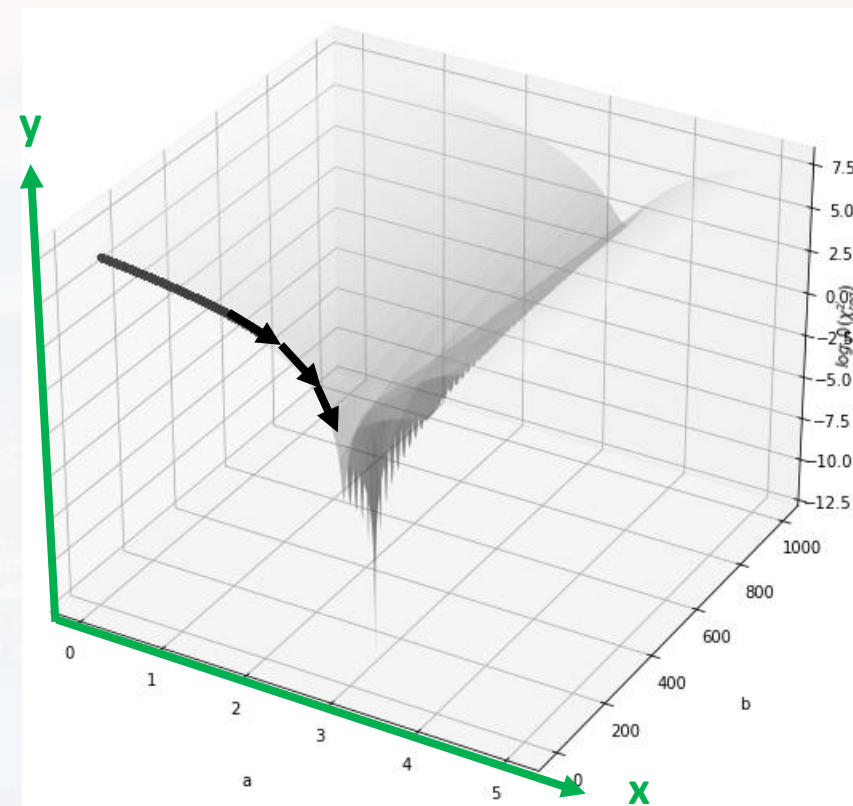$$x(t=2) = x(t=1) - \varepsilon\frac{dy}{dx}\bigg|_{x(t=1)}$$

$\varepsilon > 0$

**Vanilla**

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



$y(x(t = 1))$

$\left.\dfrac{dy}{dx}\right|_{x(t=2)}$

$x(t = 3)$

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

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

$\varepsilon > 0$

**Vanilla**

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



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

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

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

$$\varepsilon > 0$$

**Vanilla**

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



$x(t = 3)$



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

$\varepsilon > 0$

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



$y(x(t = 1))$

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

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

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

$\varepsilon > 0$

$$\left.\frac{\partial y}{\partial x_1}\right|_{x_1^*;\, x_2^*} \approx \frac{y(x_1^* + \Delta x_1,\ x_2^*) - y(x_1^* - \Delta x_1,\ x_2^*)}{2\Delta x_1}$$

$$\left.\frac{\partial y}{\partial x_2}\right|_{x_1^*;\, x_2^*} \approx \frac{y(x_1^*,\ x_2^* + \Delta x_2) - y(x_1^*,\ x_2^* - \Delta x_2)}{2\Delta x_2}$$

$$\frac{\partial y}{\partial x_1}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*} \approx \frac{y(x_1^* + \Delta x_1,\ x_2^*,...,x_N^*) - y(x_1^* - \Delta x_1,\ x_2^*,...,x_N^*)}{2\Delta x_1}$$

$$\frac{\partial y}{\partial x_2}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*} \approx \frac{y(x_1^*, x_2^* + \Delta x_2,...,x_N^*) - y(x_1^*, x_2^* - \Delta x_2,...,x_N^*)}{2\Delta x_2}$$

$$\vdots$$

$$\frac{\partial y}{\partial x_i}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*} \approx \frac{y(..., x_i^* + \Delta x_i,...,x_N^*) - y(..., x_i^* - \Delta x_i,...,x_N^*)}{2\Delta x_i}$$

$$\vdots$$

$$\frac{\partial y}{\partial x_N}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*} \approx \frac{y(x_1^*, x_2^*,...,x_N^* + \Delta x_N) - y(x_1^*, x_2^*,...,x_N^* - \Delta x_N)}{2\Delta x_N}$$

$$\begin{pmatrix} \dfrac{\partial y}{\partial x_1}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*)} \\ ... \\ \dfrac{\partial y}{\partial x_i}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*} \\ ... \\ \dfrac{\partial y}{\partial x_N}\bigg|_{x_1^*;\, x_2^*;\,...;\, x_N^*} \end{pmatrix} = grad(y)_x$$

**gradient** of $y$ wrt $x$

Outline

- The Problem

- **Gradient Descent**

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

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

$$x(t+1) = x(t) - \boldsymbol{\varepsilon}\frac{dy}{dx}\bigg|_{x(t)}$$
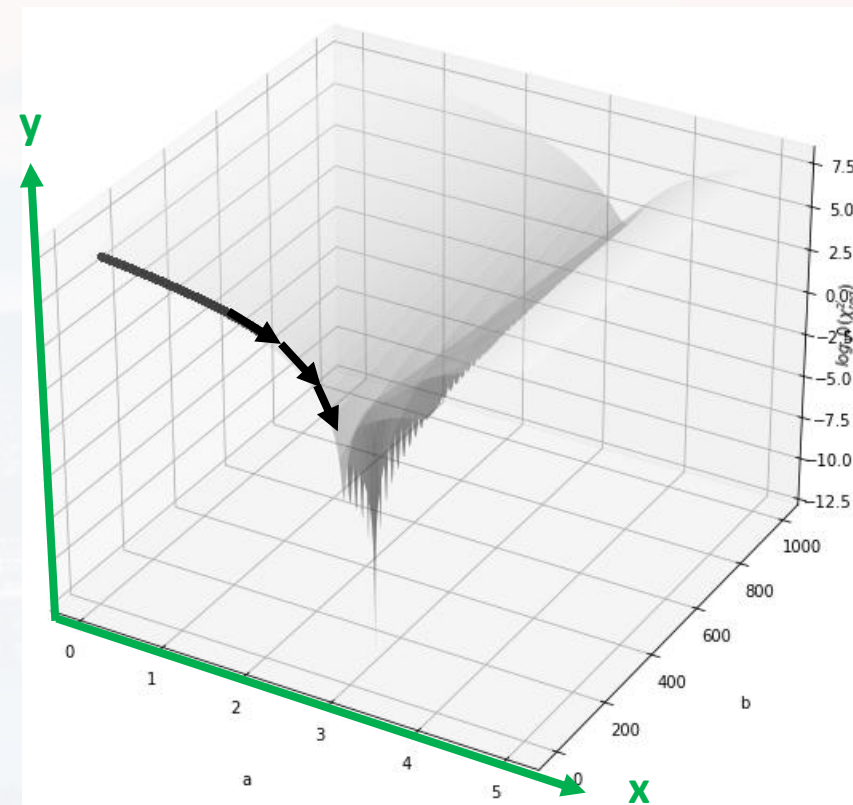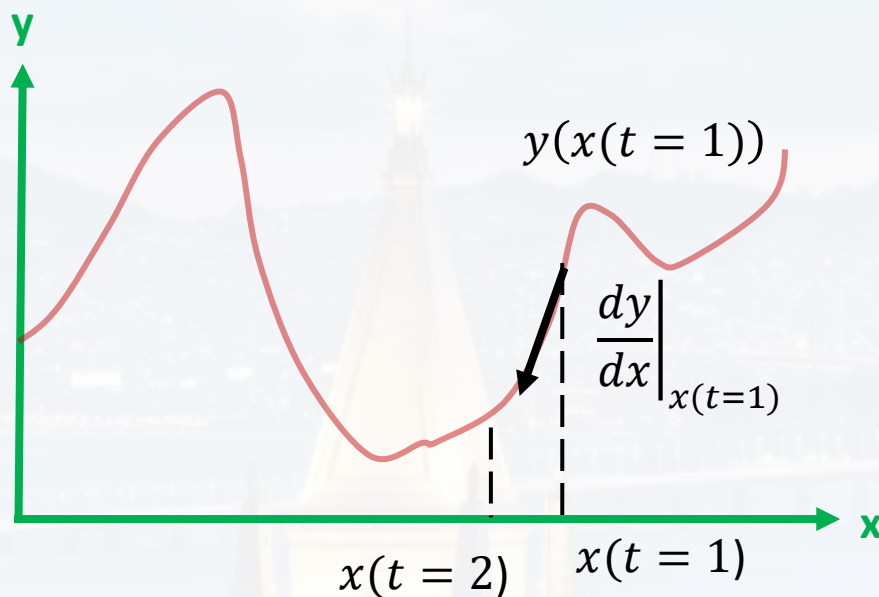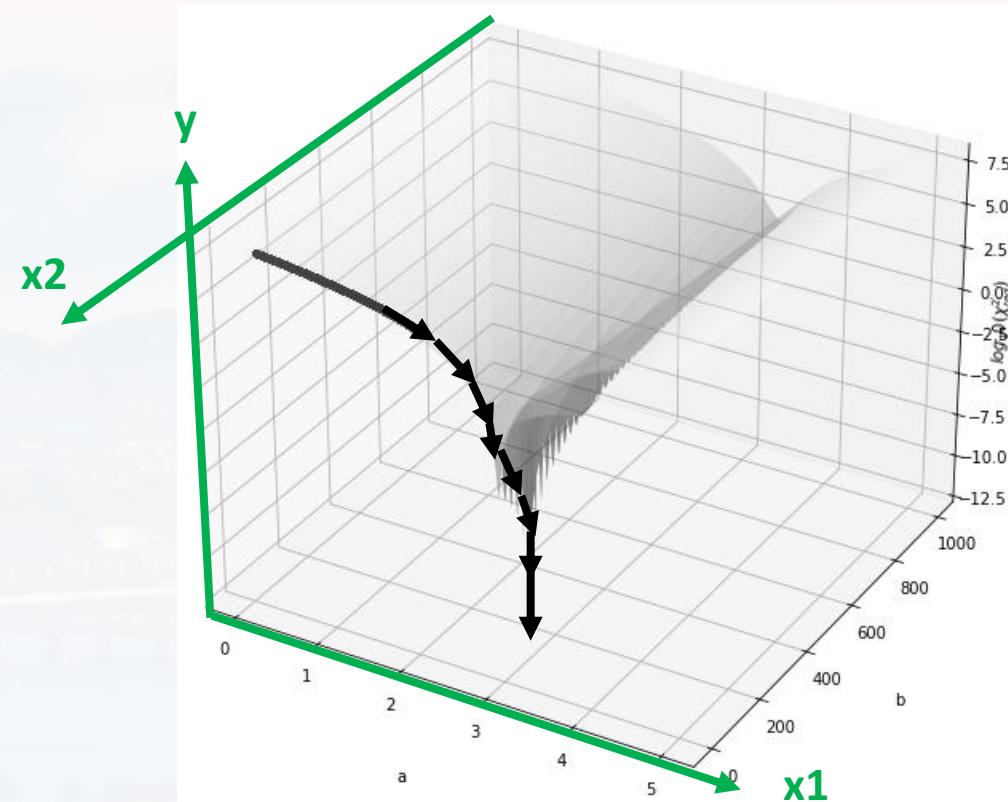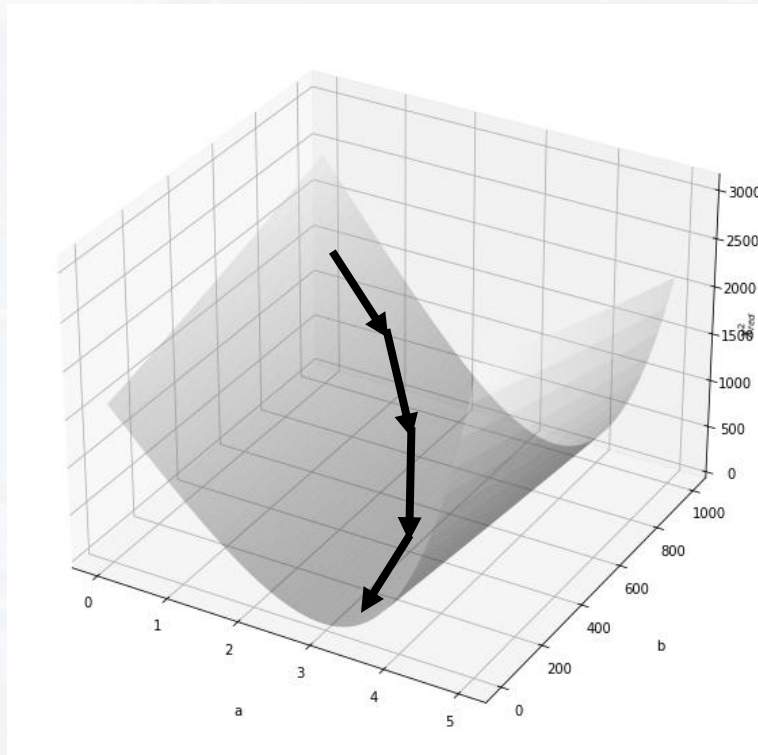
$$\boldsymbol{\varepsilon} > 0 \qquad \text{called } \textit{learning rate}$$

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

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

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

**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

**y**

**x**

$x(t = 4)$

$x(t = 3)$

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

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

**Learning Rate Schedule**

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

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

… and so on…

→ smaller $\boldsymbol{\varepsilon}$ ?



$x(t=4)$

$x(t=5)$

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

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

**Learning Rate Schedule**

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



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

defines how large
the leap $\Delta x$ is

… and so on…

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

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

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

learning rate as function of t:
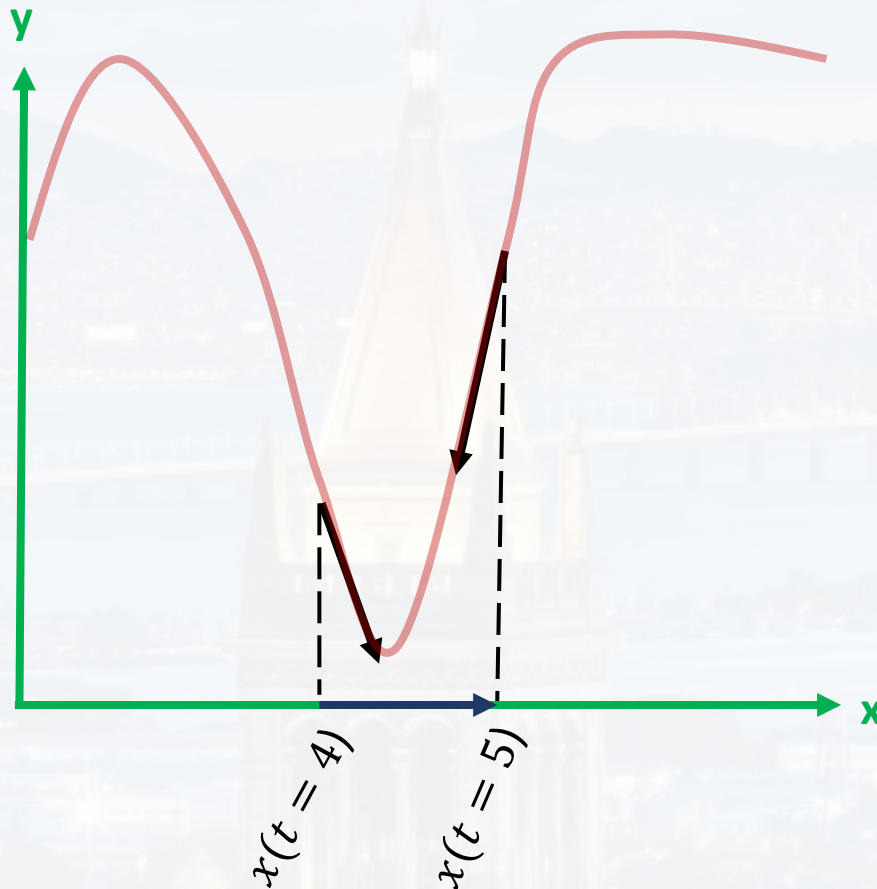
$\varepsilon > 0$     called *learning rate*

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

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

defines how large
the leap $\Delta x$ is

**Learning Rate Schedule**

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

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

learning rate as function of t:

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

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

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

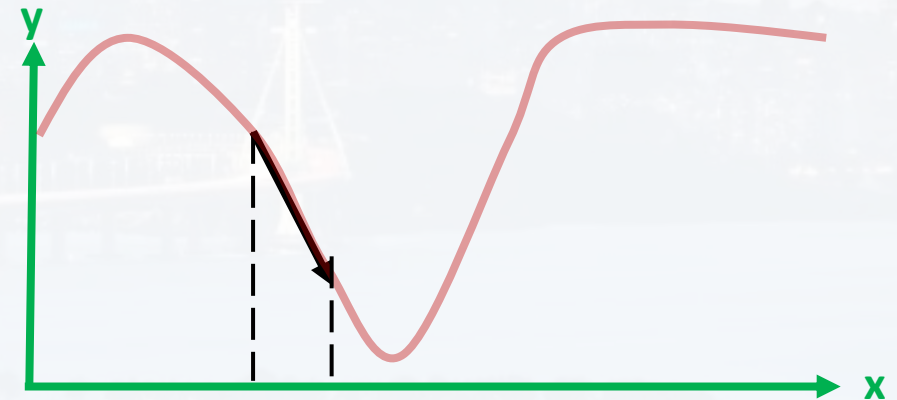can also be a stepwise function (learning rate schedule)

**Learning Rate Schedule**

<u>learning rate as function of t:</u>

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

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

can also be a stepwise function (learning rate schedule)
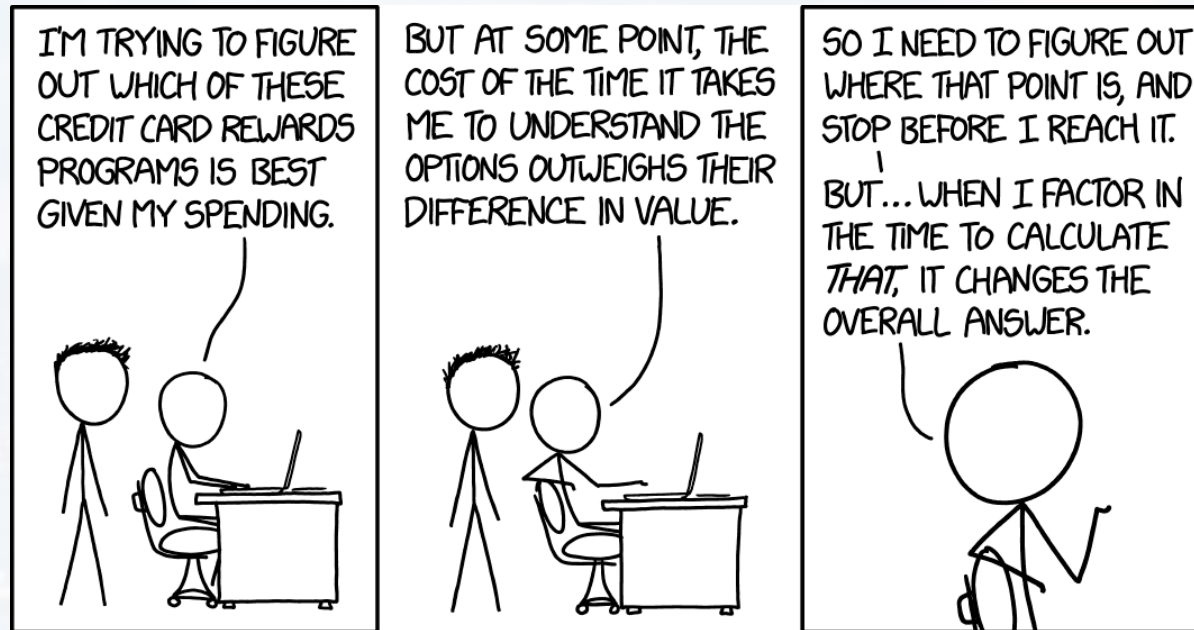


better:
adaptive gradient, aka **AdaGrad**

$$r_{t+1} = r_t + grad(y)_x \oplus grad(y)_x$$

$\oplus$:       outer product
$\delta$:       small, >0

$$\boldsymbol{\varepsilon} \rightarrow \frac{\boldsymbol{\varepsilon}}{\sqrt{r_{t+1}} + \delta}$$

Outline

- The Problem

- **Gradient Descent**

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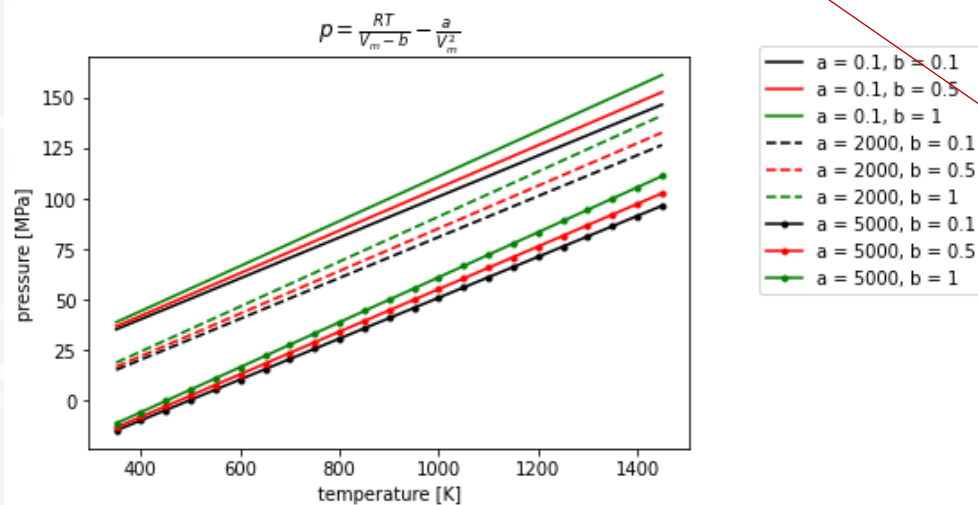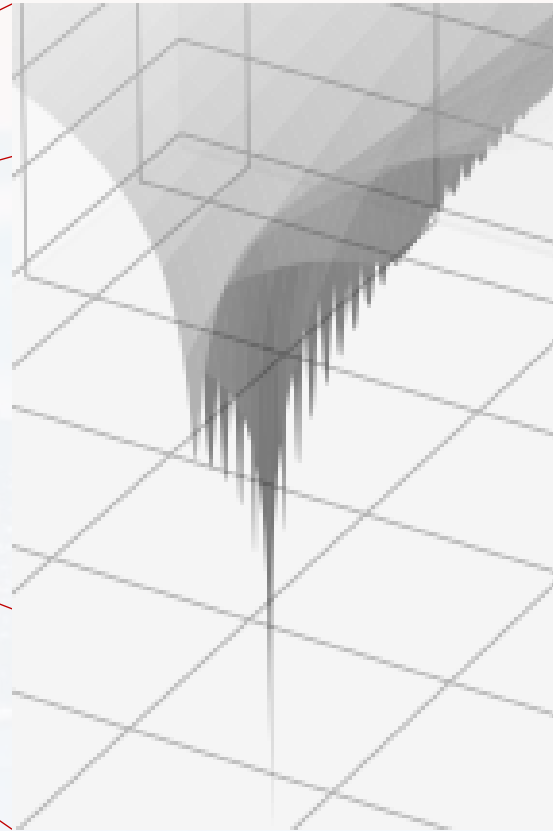  - Learning Rate Schedule
  - **Momentum**
  - L1 and L2
  - More Finetuning

**Momentum**

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

need to roll over → **momentum**

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

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 $\boldsymbol{\mu_k}$

Momentum

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 $\mu_k$

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

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

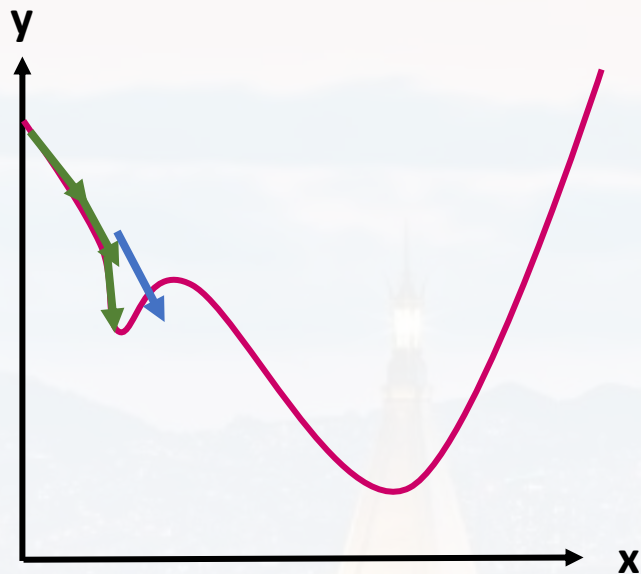$$\langle \boldsymbol{grad}(y)_{x(1)} \rangle = grad(y)_{x(1)} + \mu_0 \cdot grad(y)_{x(0)}$$

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

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

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

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

… and so on…

**weighted average** with weighting factor $\mu_k$

$\mu_0 = (0,1)$    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

L1 and L2

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

$\lambda$     *Lagrangian Multiplier*

called **L1 regularization**, or LASSO

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

**L1 and L2**

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

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

**L1 regularization**

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



| | |
|---|---|
| — | a = 0.1, b = 0.1 |
| — | a = 0.1, b = 0.5 |
| — | a = 0.1, b = 1 |
| - - - | a = 2000, b = 0.1 |
| - - - | a = 2000, b = 0.5 |
| - - - | a = 2000, b = 1 |
| —•— | a = 5000, b = 0.1 |
| —•— | a = 5000, b = 0.5 |
| —•— | a = 5000, b = 1 |

We often have even hard constrains based on the laws of physics!
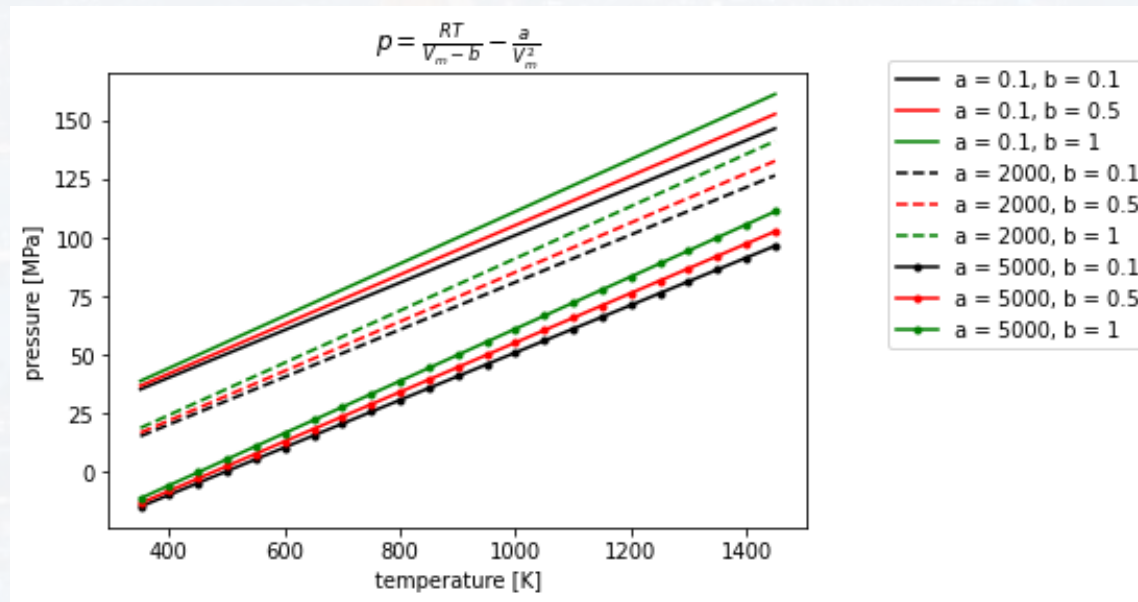
Any algorithm needs a "goal" aka **objective function** that has to be *optimized* (finding an **extreme**)

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:

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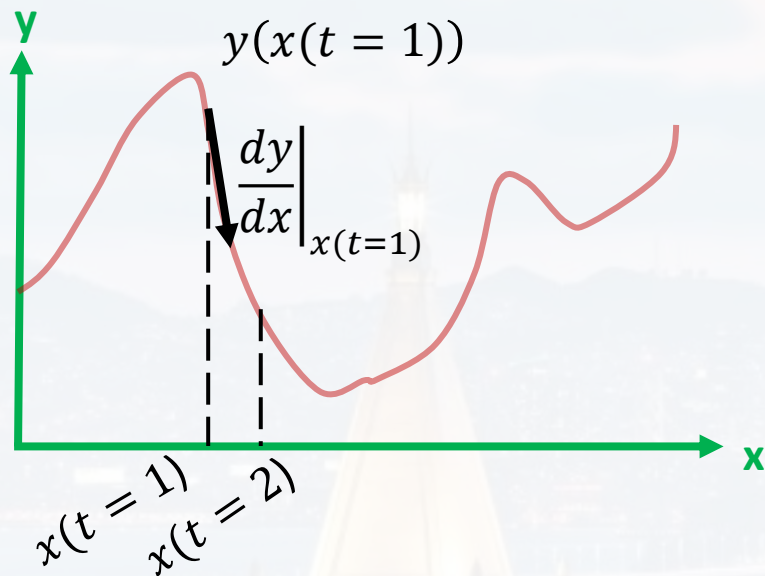
$\lambda$    *Lagrangian Multiplier*

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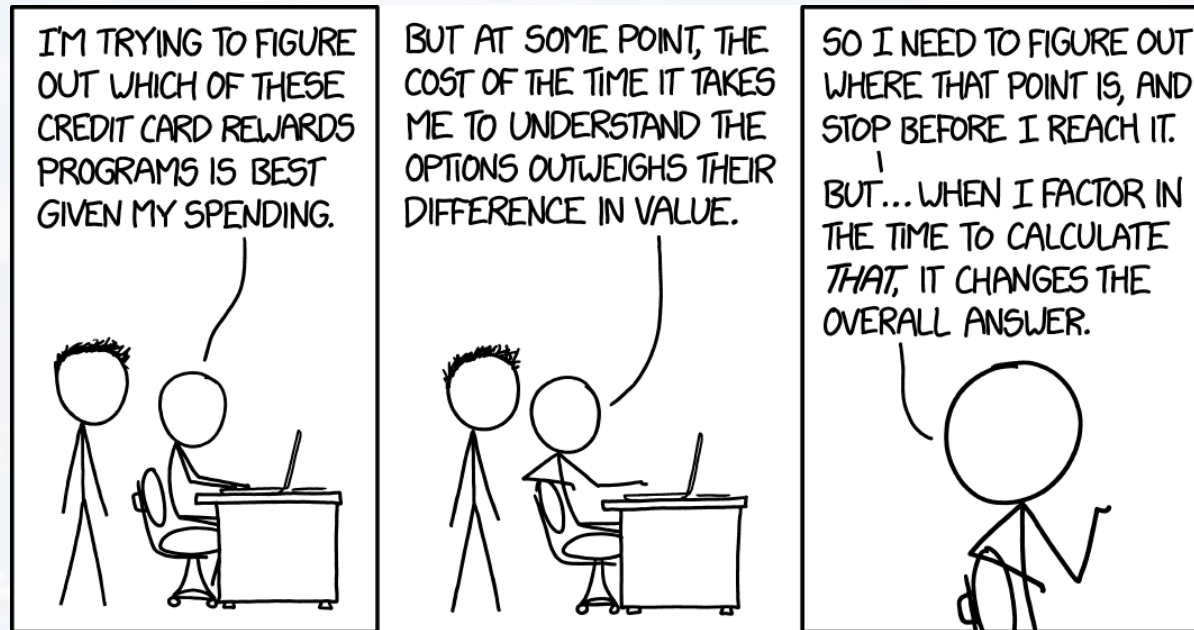
**L1 and L2 regularization**

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

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

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

- 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

低

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

**different scaling for all different directions**

momentum

$$\varepsilon \rightarrow \frac{\varepsilon}{\sqrt{r_{t+1}} + \delta}$$

adaptive gradient, aka **AdaGrad**

Multiplying 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),
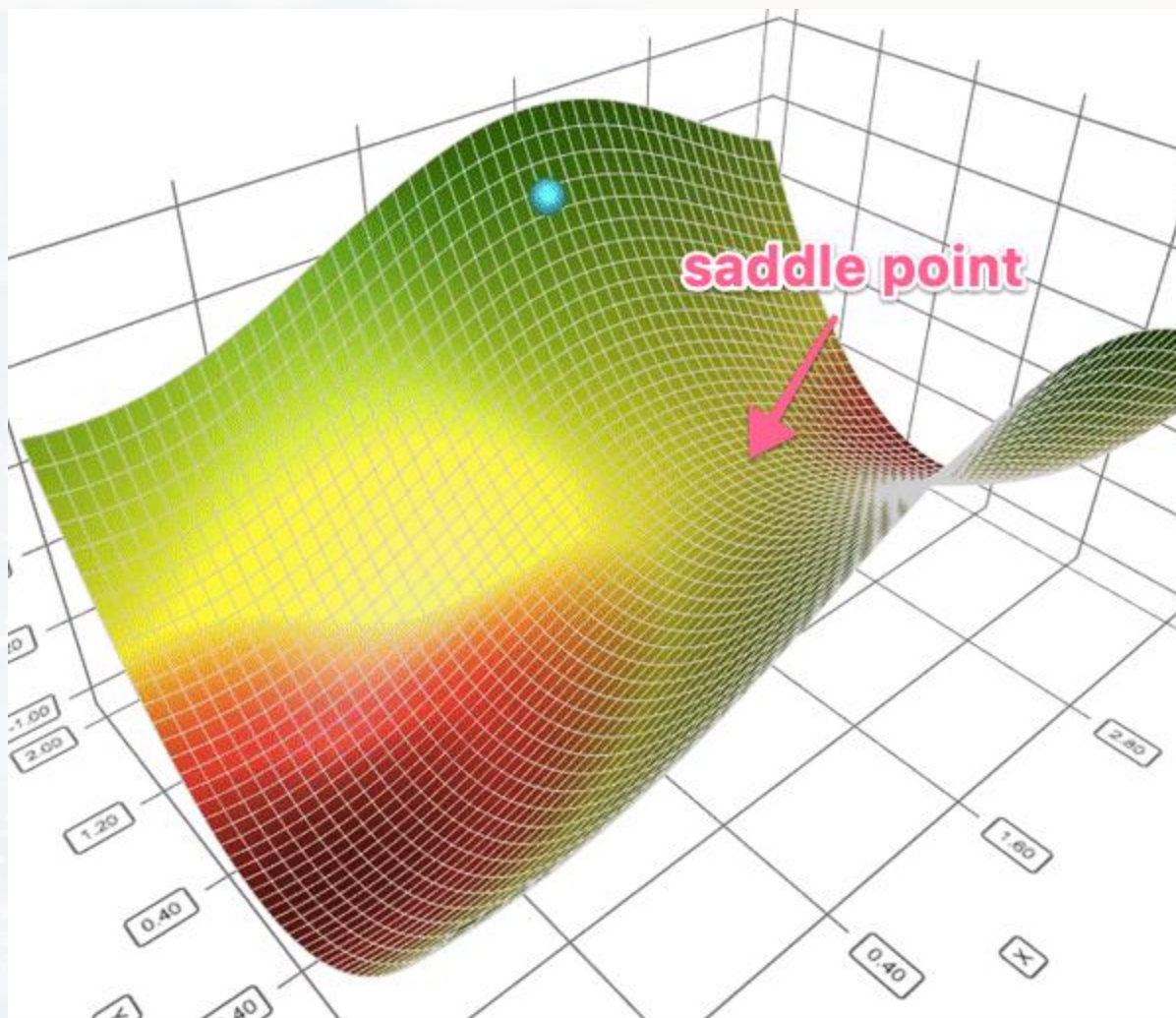AdaGrad (white),
RMSProp (green),
Adam (blue)**

TowardsDataScience

**Lili Jiang**

**More Fine Tuning**



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

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