



We want to find an extreme of an **objective function**:

- minimizing  $\chi^2 = \sum_k \frac{(\hat{y}_k - y_k)^2}{\sigma_k^2}$  curve fitting - maximizing accuracy

$\|Y - X\beta\|^2$  linear regression

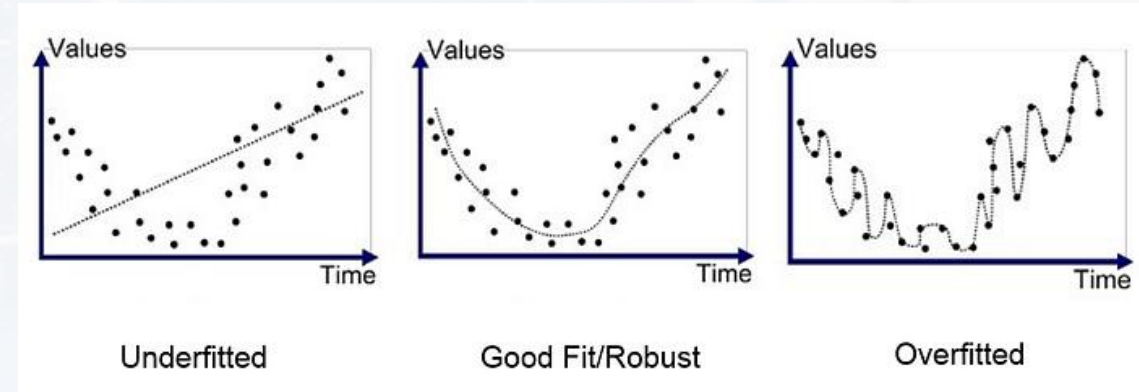
$S = - \sum_i p_i \ln p_i$  classification

$KL(p||q) = - \int p(x) \log \left[ \frac{q(x)}{\mathbf{p}(x)} \right] dx$  generation/  
encoding

We want to find an extreme of an **objective function**:

**problem:**

- algorithm might find an extreme, but for unreasonable values
- volume, mass, temperature (K) etc can only be positive
- values are extreme
- overfitting



credit: medium.com

- conservation:

$$S = - \sum_i p_i \ln p_i \quad 1 = \sum_i p_i$$

We want to find an extreme of an **objective function**:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 \}$$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda \|\beta\|^1 \}$$

the Loss Function  
 $L(X, Y, \lambda)$

L1 or **Least absolute shrinkage and selection operator**  
- encourages **sparsity** of  $\beta$   
- reduces **overfitting**

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \}$$

L2 or **Ridge**  
- **penalizes large  $\beta$**

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda \max(0, -\beta) \} \quad - \text{penalizes negative } \beta$$

...and so on

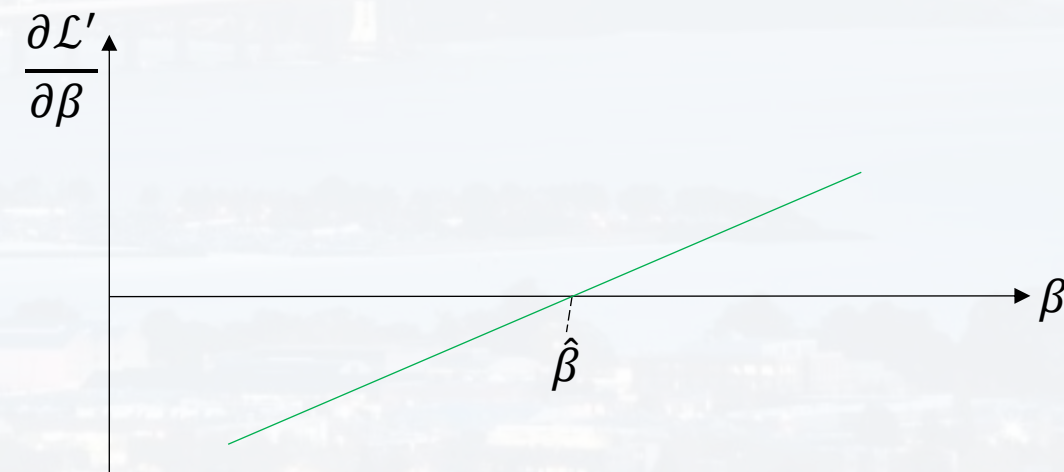
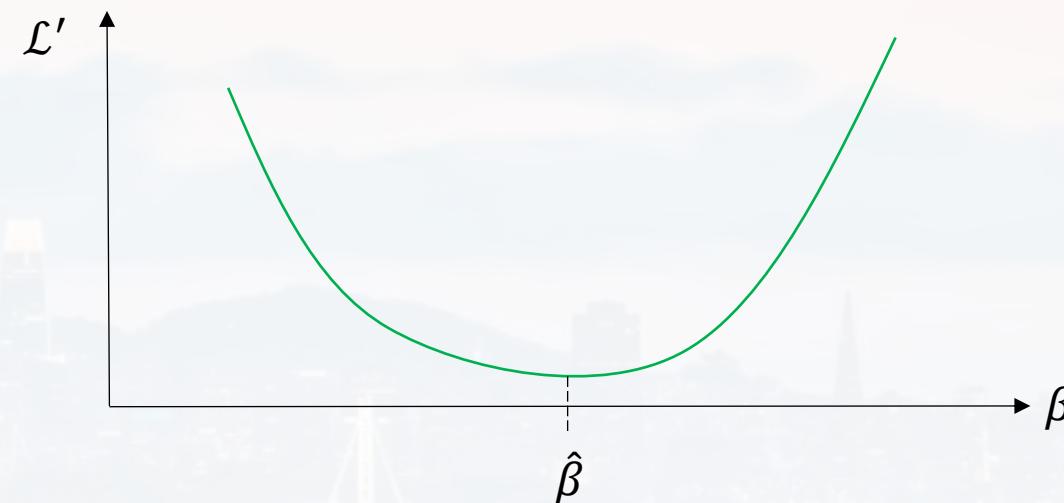


We want to find an extreme of an **objective function**:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1 + \lambda_2 \|\beta\|^2 \}$$

the Loss Function  
 $L(X, Y, \lambda)$

$$\mathcal{L}' = \|Y - X\beta\|^2$$







We want to find an extreme of an **objective function**:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1 + \lambda_2 \|\beta\|^2 \}$$

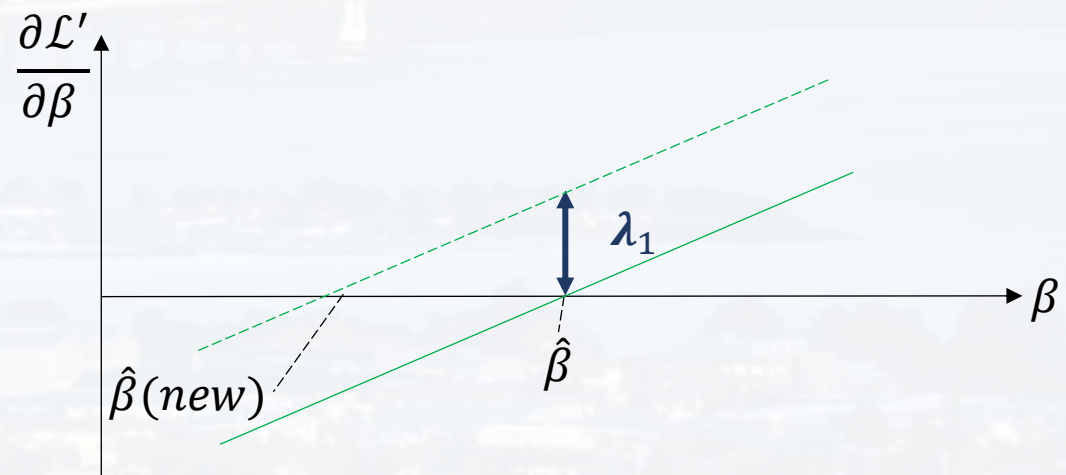
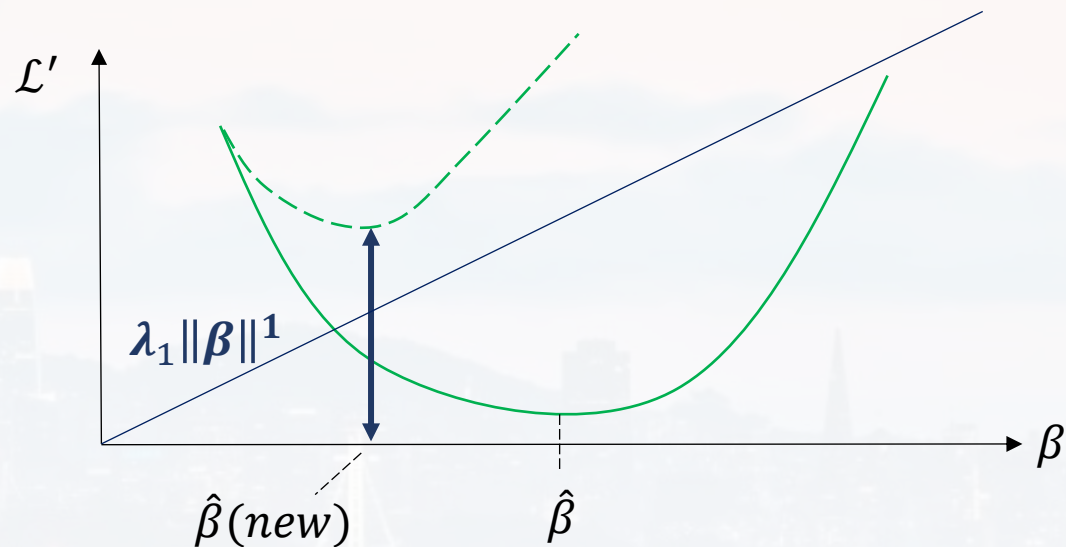
the Loss Function  
 $L(X, Y, \lambda)$

$$\mathcal{L}' = \|Y - X\beta\|^2$$

$$\begin{aligned} \frac{\partial}{\partial \beta} [\|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1] &= \frac{\partial \mathcal{L}'}{\partial \beta} + \lambda_1 \operatorname{sign}(\beta) \\ &= \frac{\partial \mathcal{L}'}{\partial \beta} \mp \lambda_1 \end{aligned}$$

shifts **if** there is a  $\beta$ , but not sensitive to its magnitude

→ large  $\lambda_1$  encourages sparsity & prevents over fitting!





We want to find an extreme of an **objective function**:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1 + \lambda_2 \|\beta\|^2 \}$$

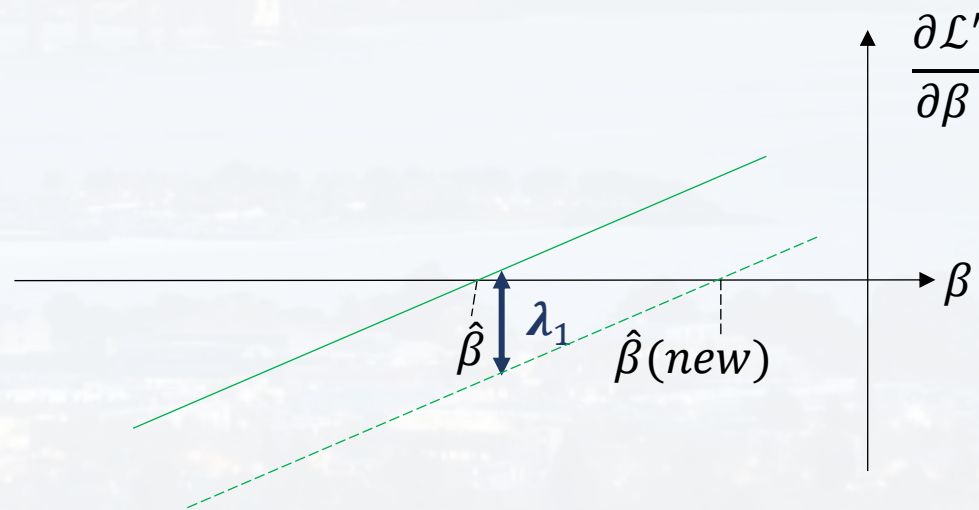
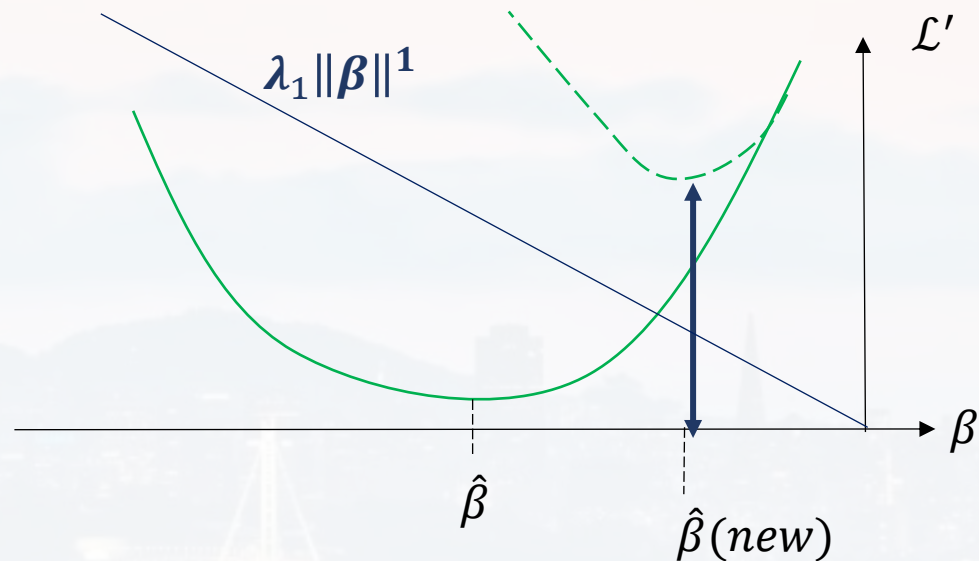
the Loss Function  
 $L(X, Y, \lambda)$

$$\mathcal{L}' = \|Y - X\beta\|^2$$

$$\begin{aligned} \frac{\partial}{\partial \beta} [\|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1] &= \frac{\partial \mathcal{L}'}{\partial \beta} + \lambda_1 \operatorname{sign}(\beta) \\ &= \frac{\partial \mathcal{L}'}{\partial \beta} \mp \lambda_1 \end{aligned}$$

shifts **if** there is a  $\beta$ , but not sensitive to its magnitude

→ large  $\lambda_1$  encourages sparsity & prevents over fitting!





We want to find an extreme of an **objective function**:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1 + \lambda_2 \|\beta\|^2 \}$$

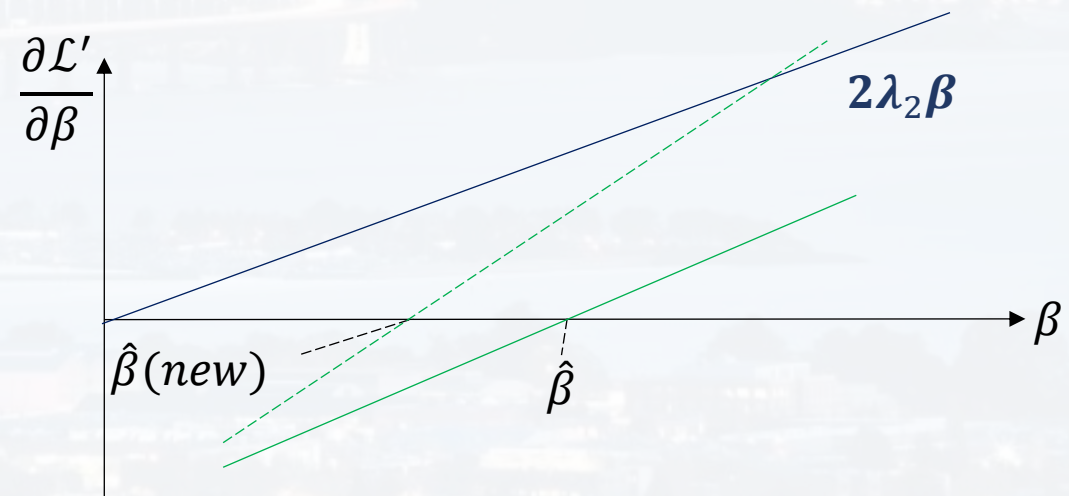
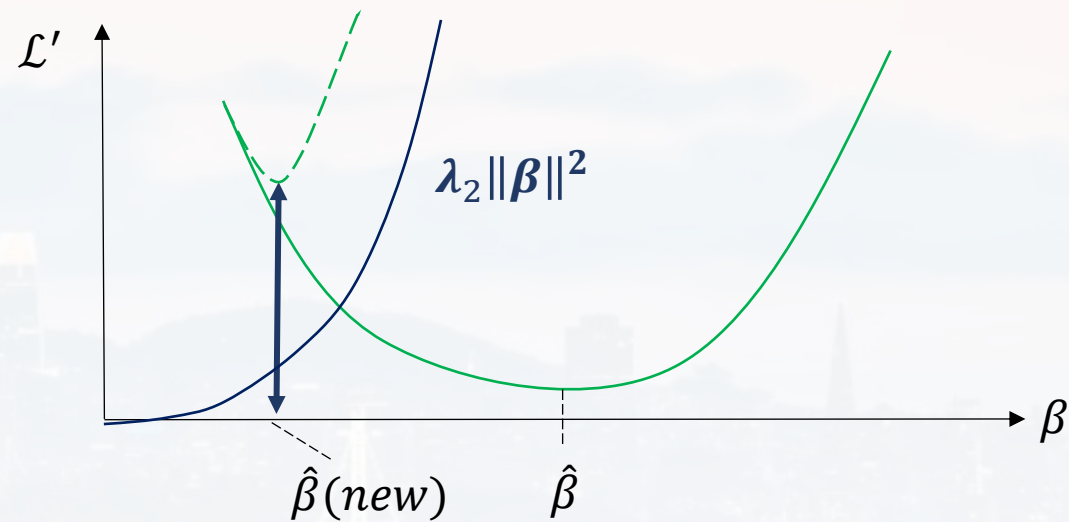
the Loss Function  
 $L(X, Y, \lambda)$

$$\mathcal{L}' = \|Y - X\beta\|^2$$

$$\frac{\partial}{\partial \beta} [\|Y - X\beta\|^2 + \lambda_2 \|\beta\|^2] = \frac{\partial \mathcal{L}'}{\partial \beta} + 2\lambda_2 \beta$$

shifts according to the  
**magnitude** of  $\beta$ ,

→ large  $\lambda_2$  encourages smaller  
magnitudes for  $\beta$







We want to find an extreme of an **objective function**:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|Y - X\beta\|^2 + \lambda_1 \|\beta\|^1 + \lambda_2 \|\beta\|^2 \}$$

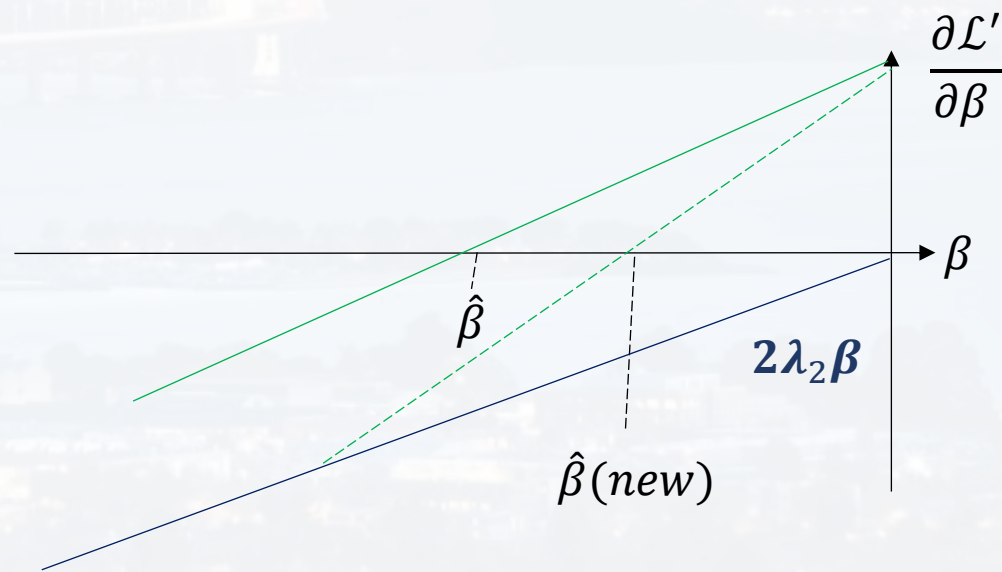
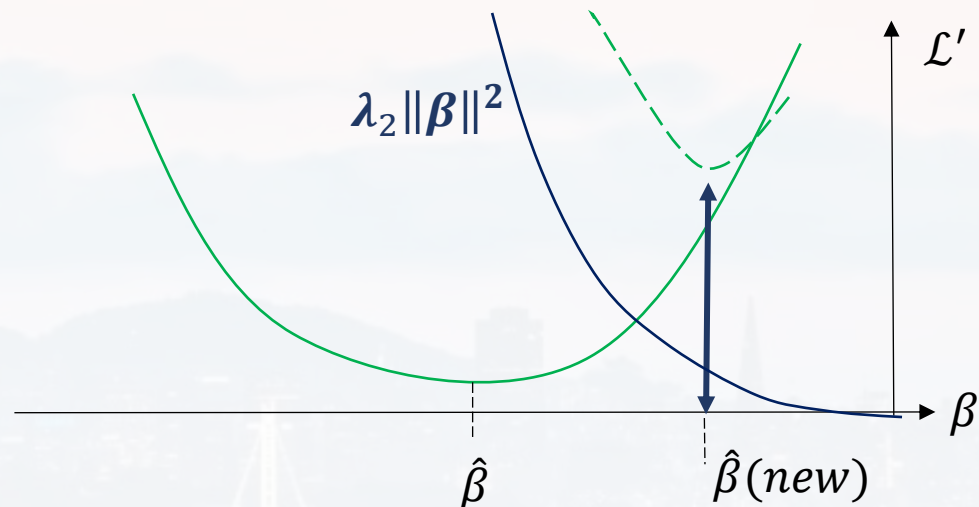
the Loss Function  
 $L(X, Y, \lambda)$

$$\mathcal{L}' = \|Y - X\beta\|^2$$

$$\frac{\partial}{\partial \beta} [\|Y - X\beta\|^2 + \lambda_2 \|\beta\|^2] = \frac{\partial \mathcal{L}'}{\partial \beta} + 2\lambda_2 \beta$$

shifts according to the  
**magnitude** of  $\beta$ ,

→ large  $\lambda_2$  encourages smaller  
magnitudes for  $\beta$

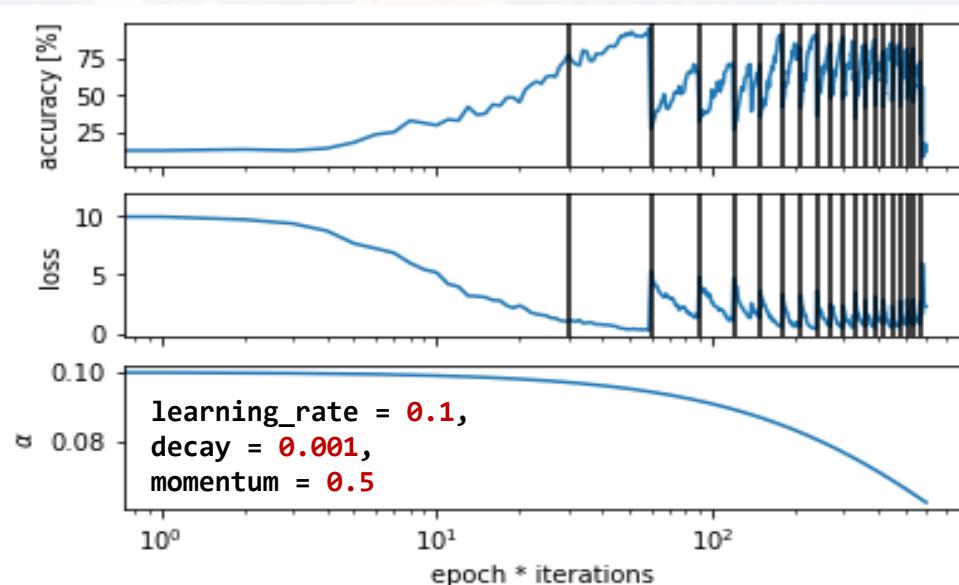




## Regularization

### note:

- how to implement L1 and L2 regularization for lin. regression in Python depends on the library (see documentation)
- “Elastic Net” balances L1 and L2:  $\lambda \left( \frac{1-\alpha}{2} \|\beta\|^2 + \alpha \|\beta\|^1 \right)$
- L1: deals with highly correlated factors (sparsity)
- L2: deals with large factors (keeps solution stable)
- regularisation is pretty common, not only for lin. models



LeNet numpy only

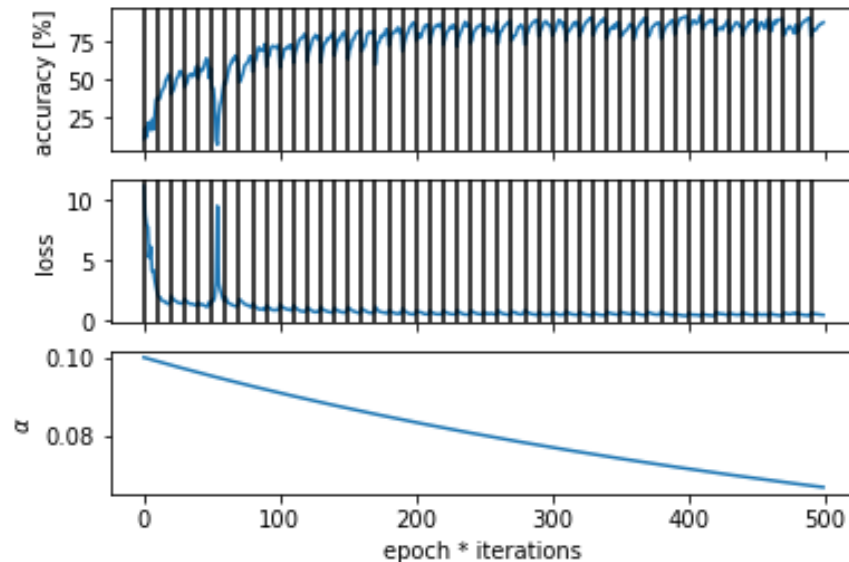
```
W = np.load('weightsC1.npy')
W[:, :, 0]
```

	0	1	2	3	4	5
0	-15923.3	-16647.4	-16277.1	-16993.1	-15715.7	-15390.1
1	9795.1	8468.91	9110.83	8205.6	9852.06	11061
2	37956.4	36572.6	37324.2	37008.4	37485.4	38515.2
3	39686.6	39131	39087.4	39614.6	39421.7	39876.7
4	25465.3	26270.1	25232.4	25836.7	25590.5	25270.9

## Regularization

**note:**  
the

- how to implement L1 and L2 regularization for lin. regression in Python depends on the library (see documentation)
- “Elastic Net” balances L1 and L2:  $\lambda \left( \frac{1-\alpha}{2} \|\beta\|^2 + \alpha \|\beta\|^1 \right)$
- L1: deals with highly correlated factors (sparsity)
- L2: deals with large factors (keeps solution stable)
- regularisation is pretty common, not only for lin. models

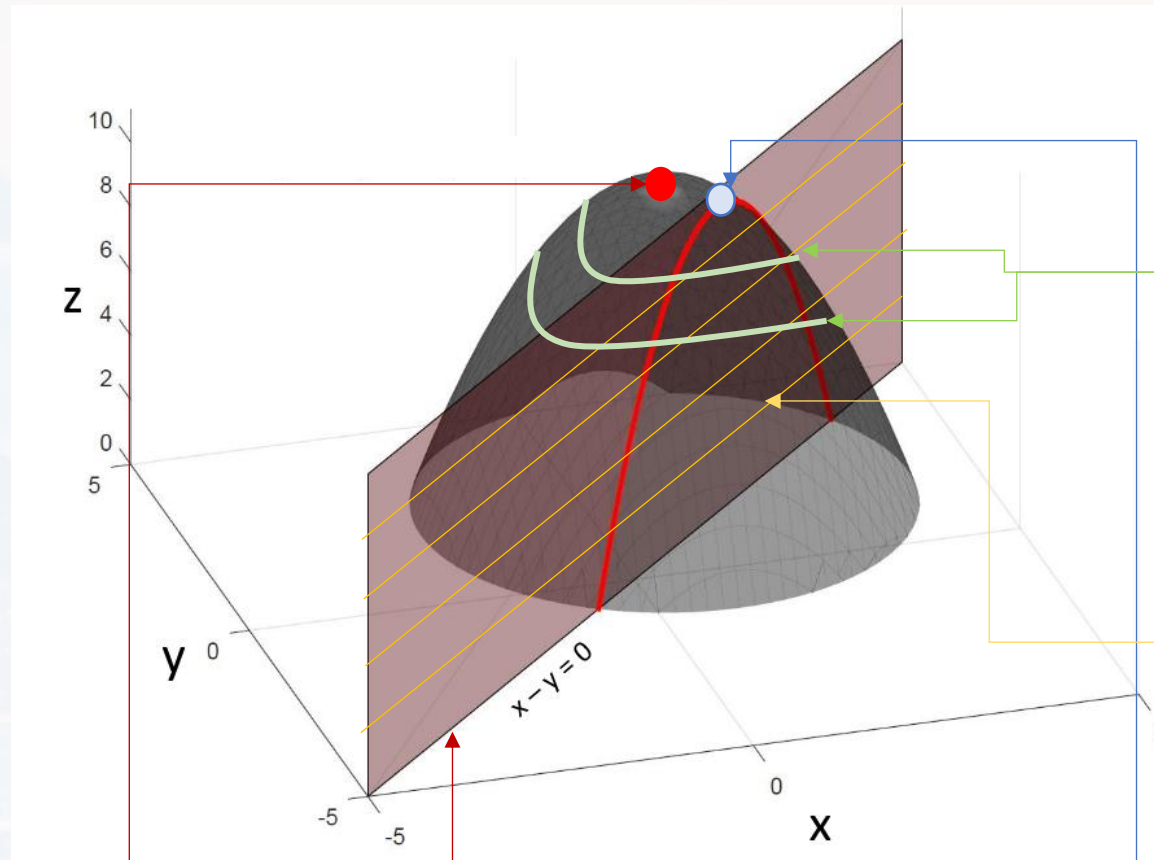


LeNet numpy only

L1 and L2 regularization

# Regularization

More about  
**Lagrangian Multiplier:**



$z = 10$  at  
 $x = 2$   
 $y = 1$

constrain  $g(x, y) = x - y = 0$

maximum of the function

**Lagrangian Multiplier**  
Examples

$$f(x, y) = z = -(x - 2)^2 - (y - 1)^2 + 10$$

level lines  $f(x, y) = \text{const}$

$$\begin{aligned} df(x, y) &= \frac{\partial f(x, y)}{\partial x} dx + \frac{\partial f(x, y)}{\partial y} dy = 0 \\ &= \text{grad} f \, d\vec{r} = 0 \end{aligned}$$

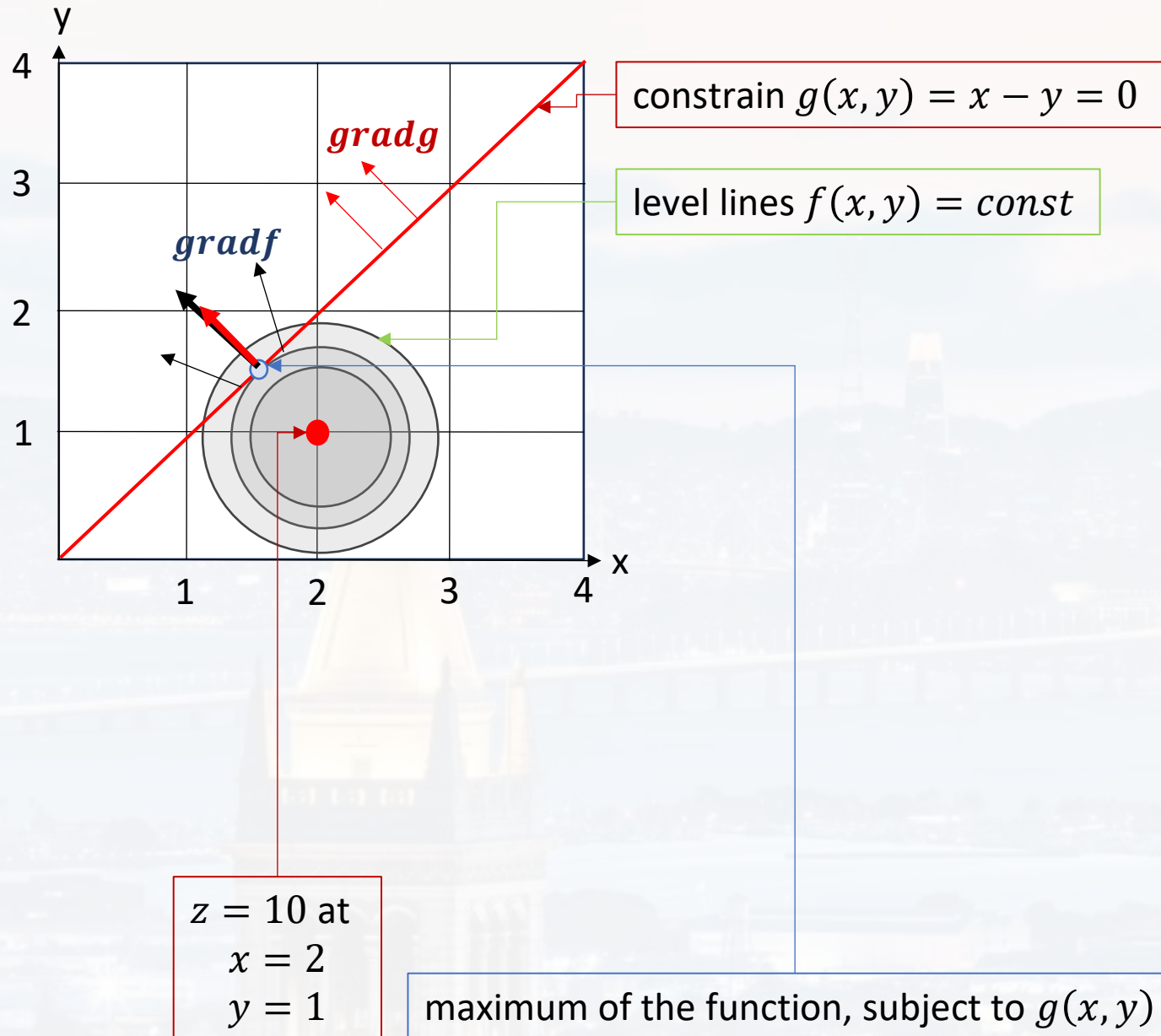
level lines  $g(x, y) = \text{const}$

$$\begin{aligned} dg(x, y) &= \frac{\partial g(x, y)}{\partial x} dx + \frac{\partial g(x, y)}{\partial y} dy = 0 \\ &= \text{grad} g \, d\vec{r} = 0 \end{aligned}$$

maximum of the function, subject to  $g(x, y)$



# Regularization



the maximum of  $f(x, y)$   
subject to  $g(x, y)$  located  
where:

$$df(x, y) = dg(x, y)$$

$$\text{grad} f \, d\vec{r} = \text{grad} g \, d\vec{r}$$

$$\text{grad} f = \text{grad} g$$

Both gradients need to point  
in the same direction  
(hence, can be multiplied with a constant,  
say  $\lambda$ )!

$$\text{grad} f = \lambda \text{grad} g$$

$\lambda$  **Lagrangian Multiplier**

the maximum of  $f(x, y)$  subject to  $g(x, y)$

$$\text{grad} f = \lambda \text{grad} g$$

$$f(x, y) - \lambda g(x, y) = \text{const}$$

**the Lagrangian**  
 $L(x, y, \lambda)$

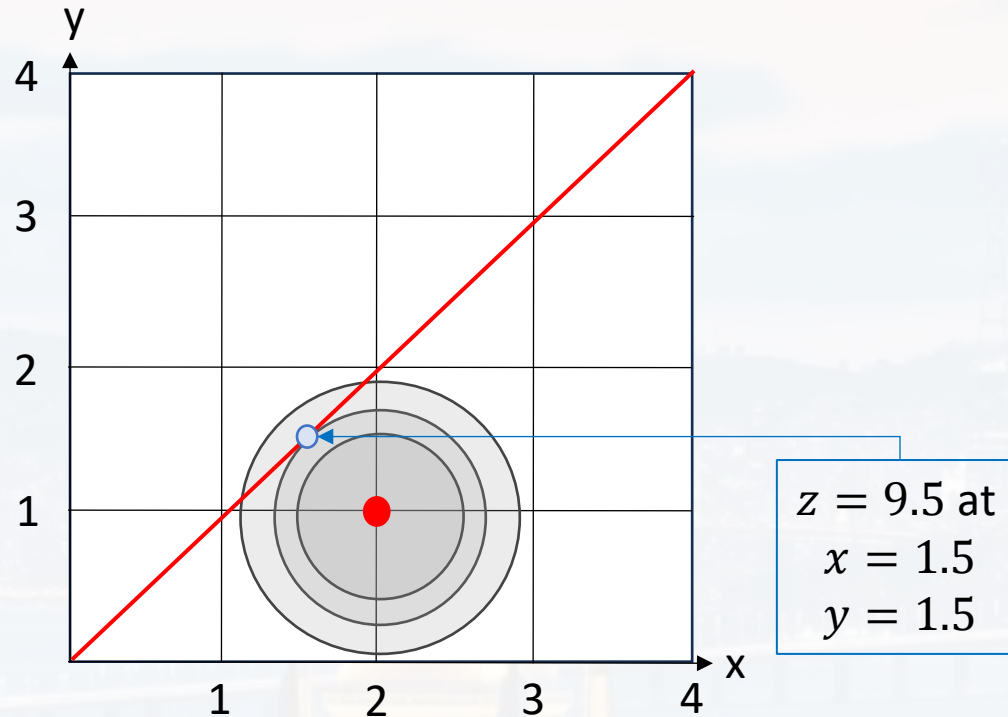
more general:

$$L(x_1, x_2, \dots, x_i, x_N, \lambda_1, \lambda_2, \dots, \lambda_k, \lambda_K) = f(x_1, x_2, \dots, x_i, x_N) - \sum_{k=1}^K \lambda_k g_k(x_1, x_2, \dots, x_i, x_N)$$

- note:**
- **$N$  dimensions and  $K \leq N$  constrains**
  - we need to solve  $N$  (from the gradient) +  $K$  equations by using the constrains
  - optimization: more robust results (most common L1 and L2 regularization, as before)
  - machine learning: including constrains in loss function (see later)

# Regularization

the maximum of  $f(x, y)$  subject to  $g(x, y)$



maximum of the function

**Lagrangian Multiplier**  
Examples

$$f(x, y) = z = -(x - 2)^2 - (y - 1)^2 + 10$$

$$\text{constrain } g(x, y) = x - y = 0$$

$$\text{constrain } x = y$$

$$x = 1.5$$

$$y = 1.5$$

$$f(1.5, 1.5) = 9.5$$

$$\text{grad} f = \lambda \text{grad} g$$

$$\frac{\partial f(x, y)}{\partial x} = \lambda \frac{\partial g(x, y)}{\partial x}$$

$$-2(x - 2) = \lambda$$

$$\frac{\partial f(x, y)}{\partial y} = \lambda \frac{\partial g(x, y)}{\partial y}$$

$$-2(y - 1) = -\lambda$$

$$y = -x + 3$$



maximum entropy of flipping a coin:

$$f(p_1, p_2) = -p_1 \ln p_1 - p_2 \ln p_2$$

subject to

$$g(p_1, p_2) = p_1 + p_2 = 1$$

absolute maximum:

Lagrangian Multiplier  
Examples

$$\frac{\partial f(p_1, p_2)}{\partial p_1} = 0$$

$$\frac{\partial f(p_1, p_2)}{\partial p_2} = 0$$

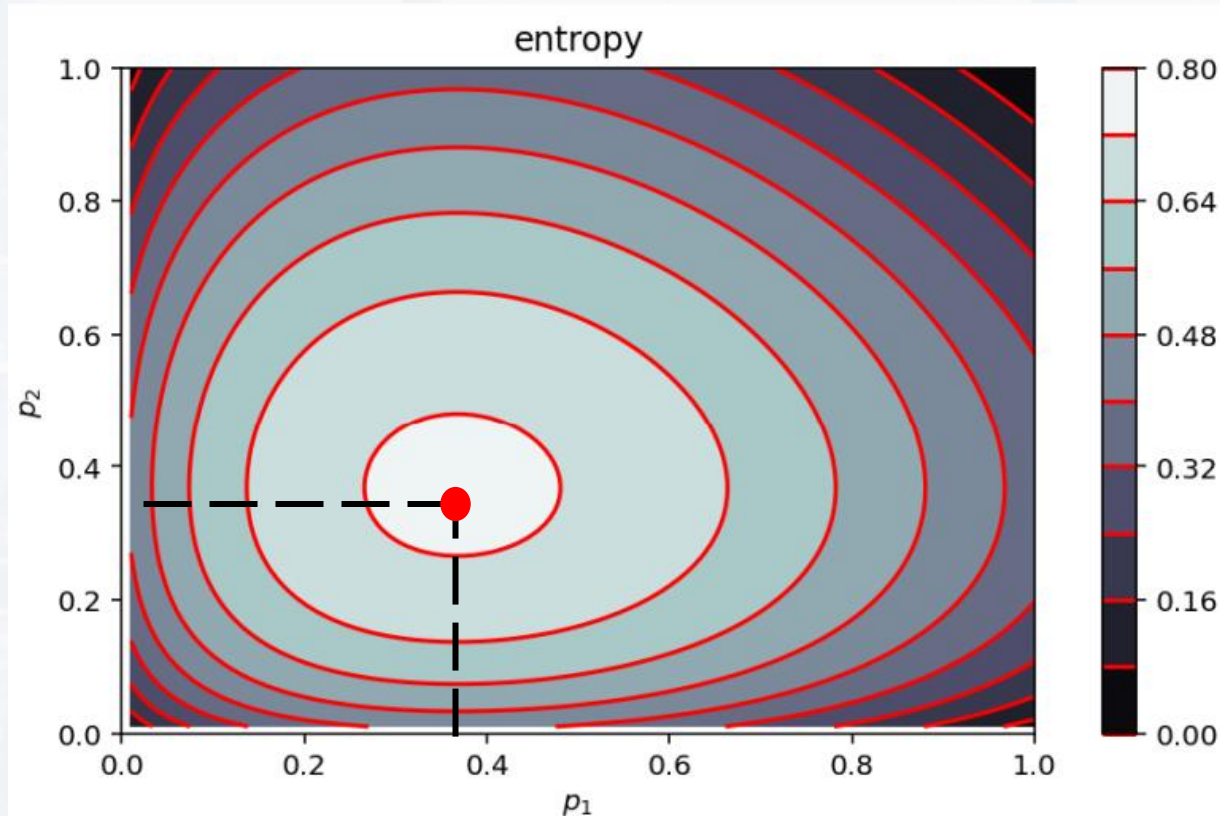
$$-\ln p_1 - 1 = 0$$

$$-\ln p_2 - 1 = 0$$

$$p_1 = p_2 = \frac{1}{e}$$

$$f\left(\frac{1}{e}, \frac{1}{e}\right) = \frac{2}{e} \approx 0.74$$

$$p_1 + p_2 = \frac{2}{e} > 1$$



maximum entropy of flipping a coin:

$$f(p_1, p_2) = -p_1 \ln p_1 - p_2 \ln p_2$$

subject to

$$g(p_1, p_2) = p_1 + p_2 = 1$$

maximum subject to  $g(p_1, p_2)$  :

$$\frac{\partial f(p_1, p_2)}{\partial p_1} = \lambda \frac{\partial g(p_1, p_2)}{\partial p_1}$$

$$\frac{\partial f(p_1, p_2)}{\partial p_2} = \lambda \frac{\partial g(p_1, p_2)}{\partial p_2}$$

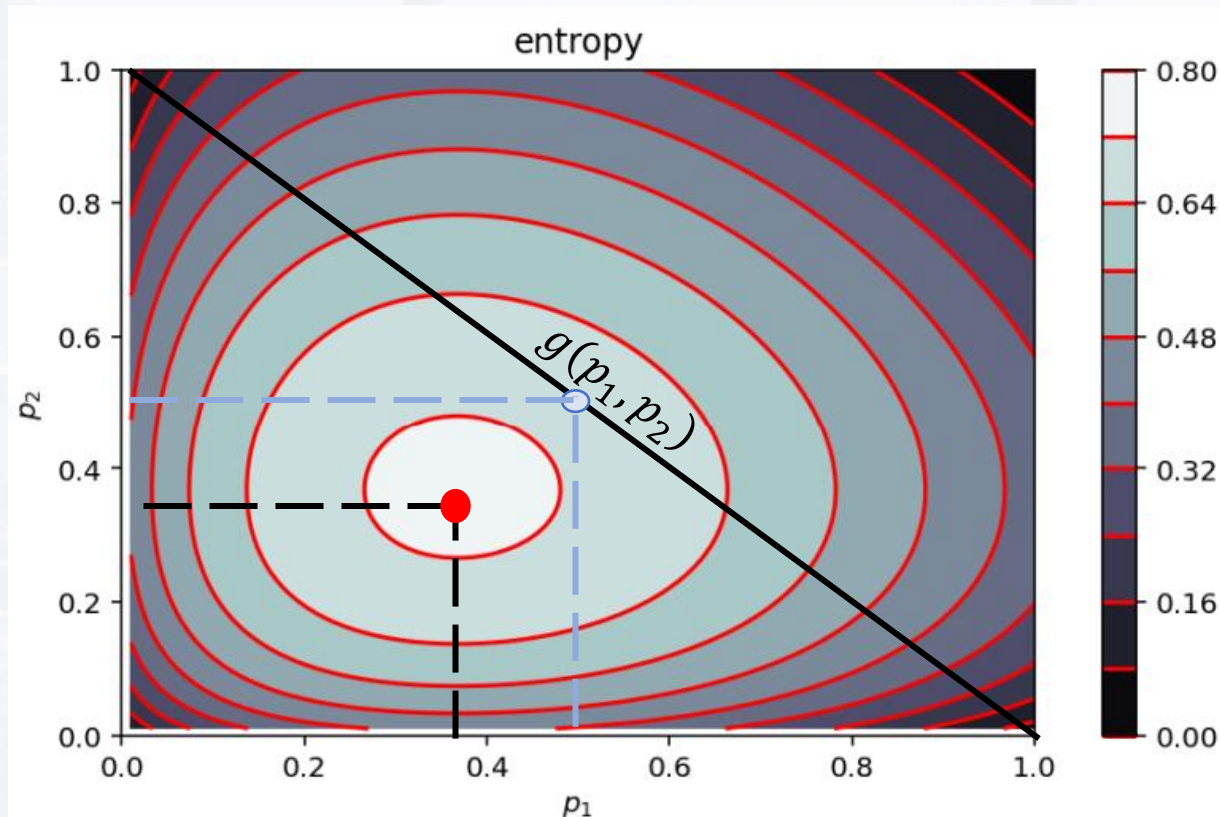
$$-\ln p_1 - 1 = \lambda \quad \quad -\ln p_2 - 1 = \lambda$$

$$p_1 = p_2$$

$$\text{constrain:} \quad p_1 + p_2 = 1$$

$$p_1 = p_2 = \frac{1}{2}$$

$$f\left(\frac{1}{2}, \frac{1}{2}\right) = \ln 2 \approx 0.69$$

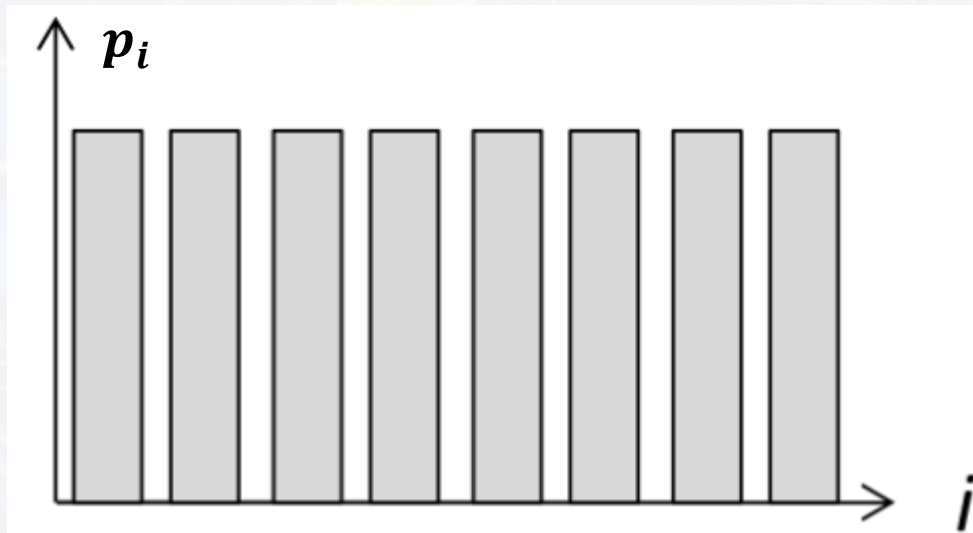


maximum entropy for  $I$  states:

$$f(p_1, \dots, p_i, \dots, p_I) = - \sum_{i=1}^I p_i \ln p_i$$

subject to

$$g(p_1, \dots, p_i, \dots, p_I) = \sum_{i=1}^I p_i = 1$$



maximum subject to  $g(p_1, \dots, p_i, \dots, p_I)$  :

$$\frac{\partial f(p_1, \dots, p_i, \dots, p_I)}{\partial p_i} = \lambda \frac{\partial g(p_1, \dots, p_i, \dots, p_I)}{\partial p_i}$$

$$-\ln p_i - 1 = \lambda \quad p_i = e^{-(1+\lambda)}$$

constrain:  $\sum_{i=1}^I e^{-(1+\lambda)} = 1$

$$I e^{-(1+\lambda)} = 1$$

$$e^{-(1+\lambda)} = \frac{1}{I}$$

**probabilities are constant!**  
**→ flat distribution!**

$$p_i = \frac{1}{I}$$



maximum entropy for  $I$  states:

$$f(p_1, \dots, p_i, \dots, p_I) = - \sum_{i=1}^I p_i \ln p_i$$

subject to

$$g_1(p_1, \dots, p_i, \dots, p_I) = \sum_{i=1}^I p_i = 1$$

if  $N$  and **total energy** is conserved

$$g_2(p_1, \dots, p_i, \dots, p_I) = \sum_{i=1}^I p_i \varepsilon_i = \frac{E_{tot}}{N} = \frac{1}{N} \sum_{i=1}^I n_i \varepsilon_i$$

$$\frac{\partial f(p_1, \dots, p_i, \dots, p_I)}{\partial p_i} = \lambda_1 \frac{\partial g_1(p_1, \dots, p_i, \dots, p_I)}{\partial p_i} + \lambda_2 \frac{\partial g_2(p_1, \dots, p_i, \dots, p_I)}{\partial p_i}$$

$$-\ln p_i - 1 = \lambda_1 + \lambda_2 \frac{\partial \sum_{j=1}^I p_j \varepsilon_j}{\partial p_i}$$

$N$ :	number of <b>indistinguishable</b> particles
$n_i$ :	number of particles in micro state $i$
$I$ :	number of states
$p_i$ :	probability of a particle being in micro state $i$
$\varepsilon_i$ :	energy in state $i$

maximum entropy for  $I$  states:

$$f(p_1, \dots, p_i, \dots, p_I) = - \sum_{i=1}^I p_i \ln p_i$$

subject to

$$g_1(p_1, \dots, p_i, \dots, p_I) = \sum_{i=1}^I p_i = 1$$

if  $N$  and **total energy** is conserved

$$g_2(p_1, \dots, p_i, \dots, p_I) = \sum_{i=1}^I p_i \varepsilon_i$$

$$-\ln p_i - 1 = \lambda_1 + \lambda_2 \frac{\partial \sum_{j=1}^I p_j \varepsilon_j}{\partial p_i}$$

$$-\ln p_i - 1 = \lambda_1 + \lambda_2 \varepsilon_i$$

$$p_i = e^{-(1+\lambda_1)} e^{-\lambda_2 \varepsilon_i}$$

from  $g_1$ :

$$p_i = \frac{1}{\sum_{i=1}^I e^{-\lambda_2 \varepsilon_i}} e^{-\lambda_2 \varepsilon_i}$$

partition function  $Z$

$$Z = \sum_{i=1}^I e^{-\lambda_2 \varepsilon_i}$$

Boltzmann distribution

$$p_i = \frac{1}{Z} e^{-\lambda_2 \varepsilon_i}$$

$N$ :	number of <b>indistinguishable</b> particles
$n_i$ :	number of particles in micro state $i$
$I$ :	number of states
$p_i$ :	probability of a particle being in micro state $i$
$\varepsilon_i$ :	energy in state $i$

## Regularization

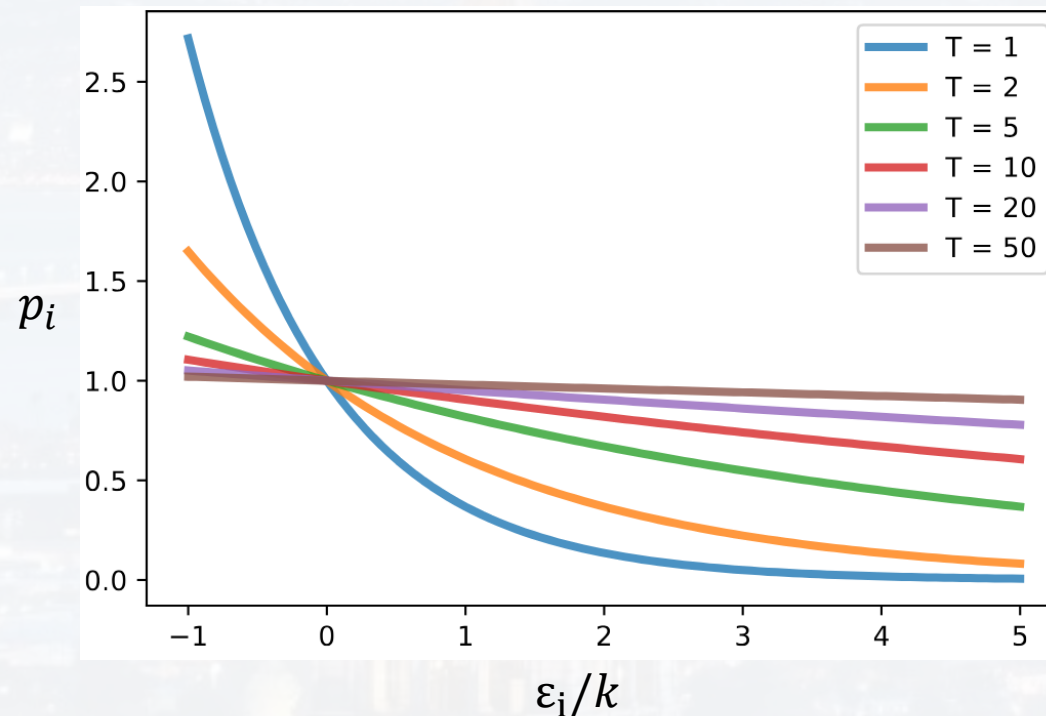
partition function  $Z$

$$Z = \sum_{i=1}^I e^{-\lambda_2 \varepsilon_i}$$

Boltzmann distribution

$$p_i = \frac{1}{Z} e^{-\lambda_2 \varepsilon_i}$$

one can show:  $\lambda_2 = \frac{1}{kT}$



### Lagrangian Multiplier Examples

$N$ :	number of <b>indistinguishable</b> particles
$n_i$ :	number of particles in micro state $i$
$I$ :	number of states
$p_i$ :	probability of a particle being in micro state $i$
$\varepsilon_i$ :	energy in state $i$

- note:**
- for  $T \rightarrow \infty$ ,  $Z \rightarrow I$ , i. e. higher states become more accessible and  $p_i \rightarrow \frac{1}{I}$
  - we used maximum entropy: **equilibrium** state for large  $N$
  - $N$  and  $E_{\text{tot}}$  are constant
  - ANN: **softmax layer** for classification probabilities (see later)



# Regularization

