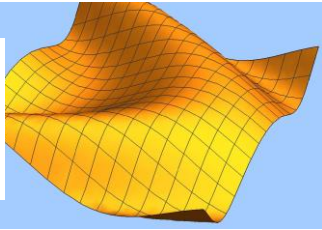


Inheriting regularization through Knowledge Distillation



Statistical
Physics of
Computation
Laboratory



Luca Saglietti, Lenka Zdeborová
SPOC lab – EPFL

MSLS 2021

EPFL

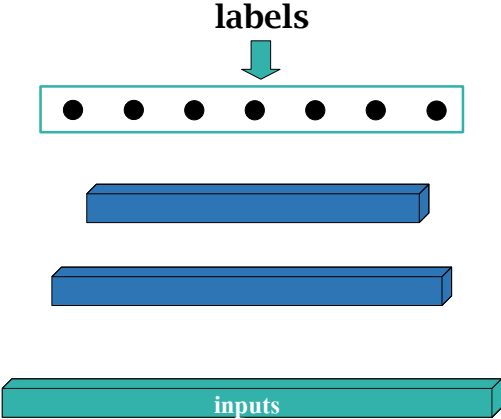
Knowledge distillation for network compression



G. Hinton, O. Vinyals, J. Dean, NIPS 2015

R. Anil, G. Hinton et al., ICLR 2018

Knowledge distillation for network compression



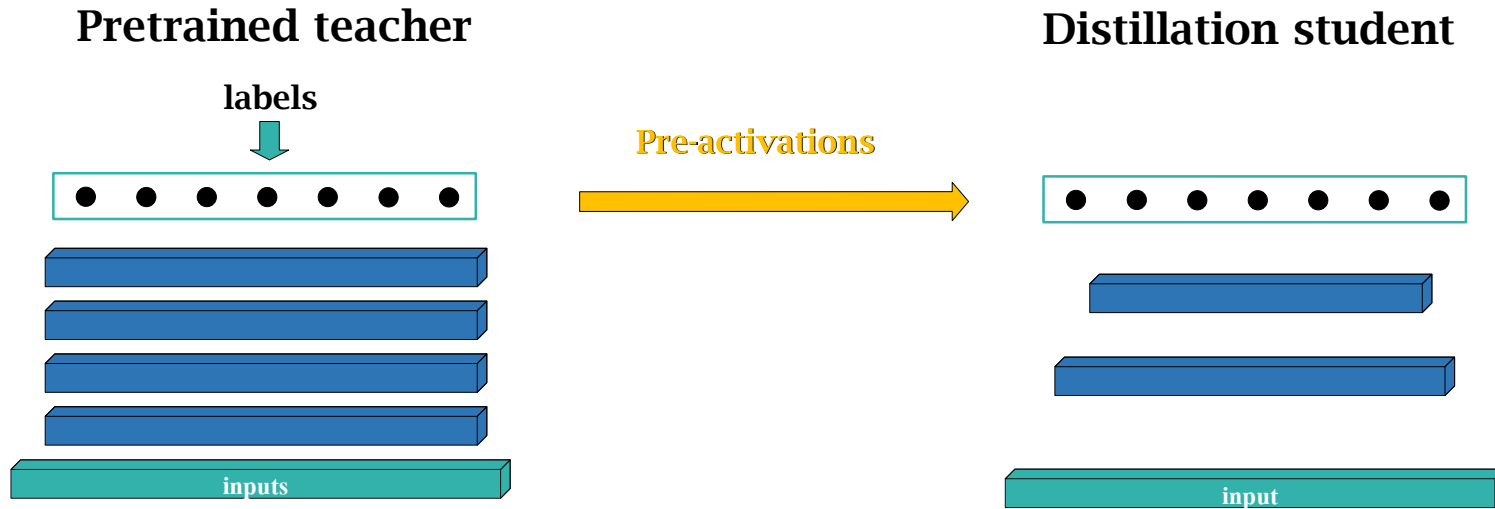
Small NN → optimization harder → **POOR GENERALIZATION**

Knowledge distillation for network compression



Larger NN → optimization bias (**implicit regularization**) → **GOOD GENERALIZATION**

Knowledge distillation for network compression



Larger NN → optimization bias (**implicit regularization**) → **GOOD GENERALIZATION**

The **pre-activations** retain:

- **uncertainty** estimation
- **relational information** between categories
- **reweight** the training samples

KD loss function

$$\mathcal{L}_{KD} = \sum_{\mu} \left((1 - \chi) \mathcal{H}(y^{\mu}, f(\mathbf{x}; \mathbf{W})) + \chi \mathcal{H}(\tilde{f}(\mathbf{x}; \tilde{\mathbf{W}}, T), f(\mathbf{x}; \mathbf{W}, T)) \right) + \lambda \|\mathbf{W}\|$$

KD loss function


$$\mathcal{L}_{KD} = \sum_{\mu} \left((1 - \chi) \overset{\text{usual logistic regression with true labels}}{\mathcal{H}(y^{\mu}, f(\mathbf{x}; \mathbf{W}))} + \chi \overset{\text{logistic regression with teacher pre-activations}}{\mathcal{H}(\tilde{f}(\mathbf{x}; \tilde{\mathbf{W}}, T), f(\mathbf{x}; \mathbf{W}, T))} \right) + \lambda \|\mathbf{W}\|$$

■ **Cross-entropy:** $\mathcal{H}(y, p) = - (y \log(p) + (1 - y) \log(1 - p))$

KD loss function

$$\mathcal{L}_{KD} = \sum_{\mu} \left((1 - \chi) \mathcal{H}(y^{\mu}, f(\mathbf{x}; \mathbf{W})) + \chi \mathcal{H}(\tilde{f}(\mathbf{x}; \tilde{\mathbf{W}}, T), f(\mathbf{x}; \mathbf{W}, T)) \right) + \lambda \|\mathbf{W}\|$$

usual logistic regression with true labels logistic regression with teacher pre-activations


 **Cross-entropy:** $\mathcal{H}(y, p) = -(y \log(p) + (1 - y) \log(1 - p))$

 **KD mixing parameter**

KD loss function

$$\mathcal{L}_{KD} = \sum_{\mu} \left((1 - \chi) \mathcal{H}(y^{\mu}, f(\mathbf{x}; \mathbf{W})) + \chi \mathcal{H}(\tilde{f}(\mathbf{x}; \tilde{\mathbf{W}}, T), f(\mathbf{x}; \mathbf{W}, T)) \right) + \lambda \|\mathbf{W}\|$$

usual logistic regression with true labels logistic regression with teacher pre-activations

 **Cross-entropy:** $\mathcal{H}(y, p) = -(y \log(p) + (1 - y) \log(1 - p))$

 **KD mixing parameter**

 **Direct (explicit) regularization**

Theory?

Studying **distillation**: stat phys approach

⇒ **2-level problem:**

a) pre-train teacher \tilde{W}

b) train student W

⇒ Step **a)** is **unaffected** by step **b)**

⇒ Both levels share the **training set**

Studying **distillation**: stat phys approach

⇒ **2-level problem:**

a) pre-train teacher \tilde{W}

b) train student W

⇒ Step **a)** is **unaffected** by step **b)**

⇒ Both levels share the **training set**

⇒ **Franz-Parisi potential formalism:**

S. Franz and G. Parisi, PRL 1997

$$S_{FP} = \int d\tilde{W} \frac{e^{-\tilde{\beta} \tilde{E}(\tilde{W})}}{\tilde{Z}(\tilde{\beta})} \log \int d\mathbf{W} e^{-\beta E(\mathbf{W}, \tilde{W})}$$

$$E(\mathbf{W}, \xi = \{\mathbf{x}^\mu, y^\mu\}) = \sum_{\mu} \ell(\hat{y}(\mathbf{W}, \mathbf{x}^\mu), y^\mu) + \lambda \|\mathbf{W}\|$$

Studying **distillation**: stat phys approach

⇒ **2-level problem:**

a) pre-train teacher \tilde{W}

b) train student W

⇒ Step **a)** is **unaffected** by step **b)**

⇒ Both levels share the **training set**

⇒ **Franz-Parisi potential formalism:**

$$\mathcal{S}_{FP} = \int d\tilde{W} \frac{e^{-\tilde{\beta} \tilde{E}(\tilde{W})}}{\tilde{Z}(\tilde{\beta})} \log \int dW e^{-\beta E(W, \tilde{W})}$$

S. Franz and G. Parisi, PRL 1997

Studying **distillation**: stat phys approach

⇒ **2-level problem**:

a) pre-train teacher \tilde{W}

b) train student W

⇒ Step **a)** is **unaffected** by step **b)**

⇒ Both levels share the **training set**

⇒ **Franz-Parisi potential formalism**:

$$S_{FP} = \int d\tilde{W} \frac{e^{-\tilde{\beta} \tilde{E}(\tilde{W})}}{\tilde{Z}(\tilde{\beta})} \log \int dW e^{-\beta E(W, \tilde{W})}$$

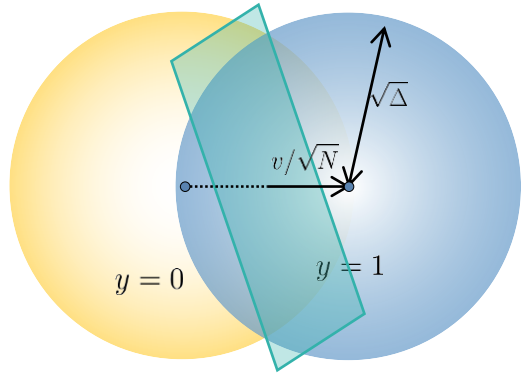
→ **Quenched** and **annealed** disorder → **REPLICA METHOD**

S. Franz and G. Parisi, PRL 1997

Model assumptions

Data model: Isotropic Gaussian mixture
(2 clusters, M points in dimension N)

Learning model: L2-regularized logistic regression

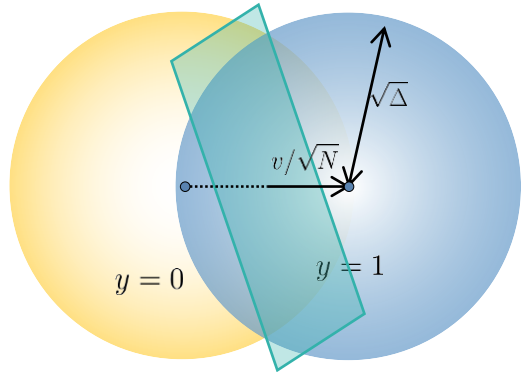


Model assumptions

Data model: Isotropic Gaussian mixture
(2 clusters, M points in dimension N)

Learning model: L2-regularized logistic regression

Asymptotic limit: $N, M \rightarrow \infty$ $M/N = \alpha$



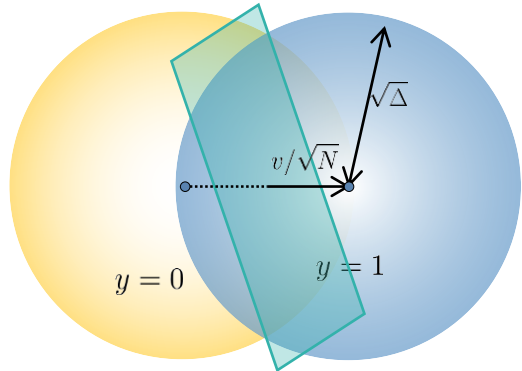
Tuning regularization intensity λ is key!!!

Model assumptions

Data model: Isotropic Gaussian mixture
(2 clusters, M points in dimension N)

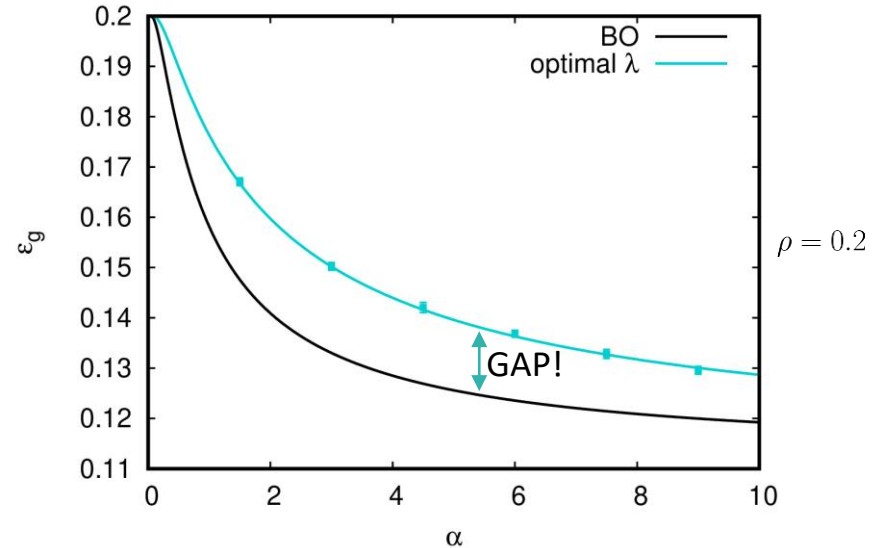
Learning model: L2-regularized logistic regression

Asymptotic limit: $N, M \rightarrow \infty$ $M/N = \alpha$



Tuning **regularization intensity λ** is key!!!

Unbalanced clusters: $y^\mu = \text{Bern}(\rho)$ **hard!**

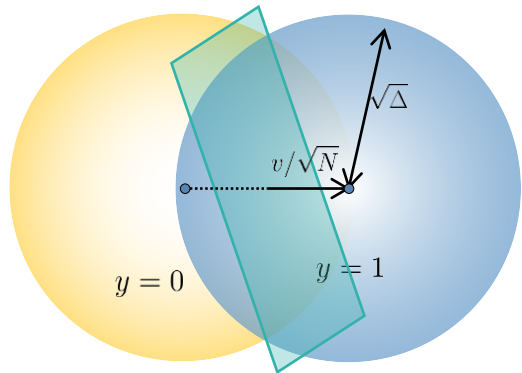


Model assumptions

Data model: Isotropic Gaussian mixture
(2 clusters, M points in dimension N)

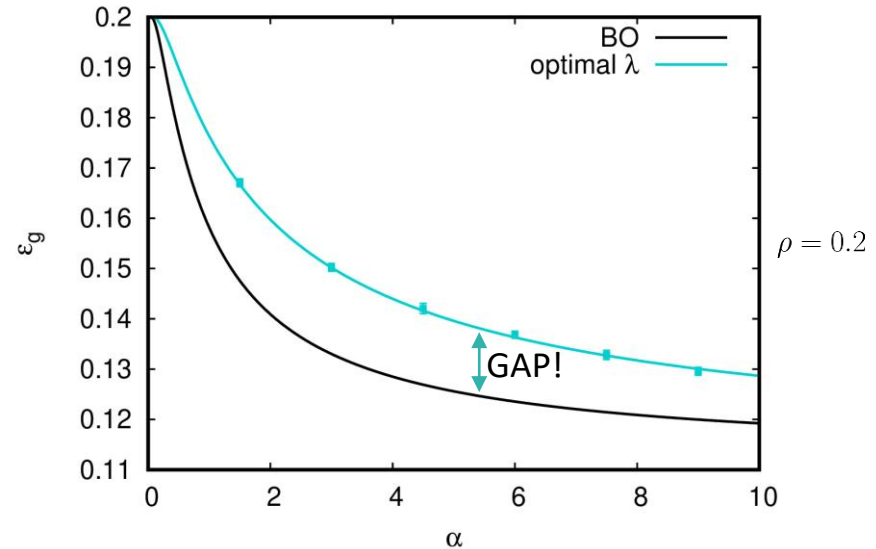
Learning model: L2-regularized logistic regression

Asymptotic limit: $N, M \rightarrow \infty$ $M/N = \alpha$



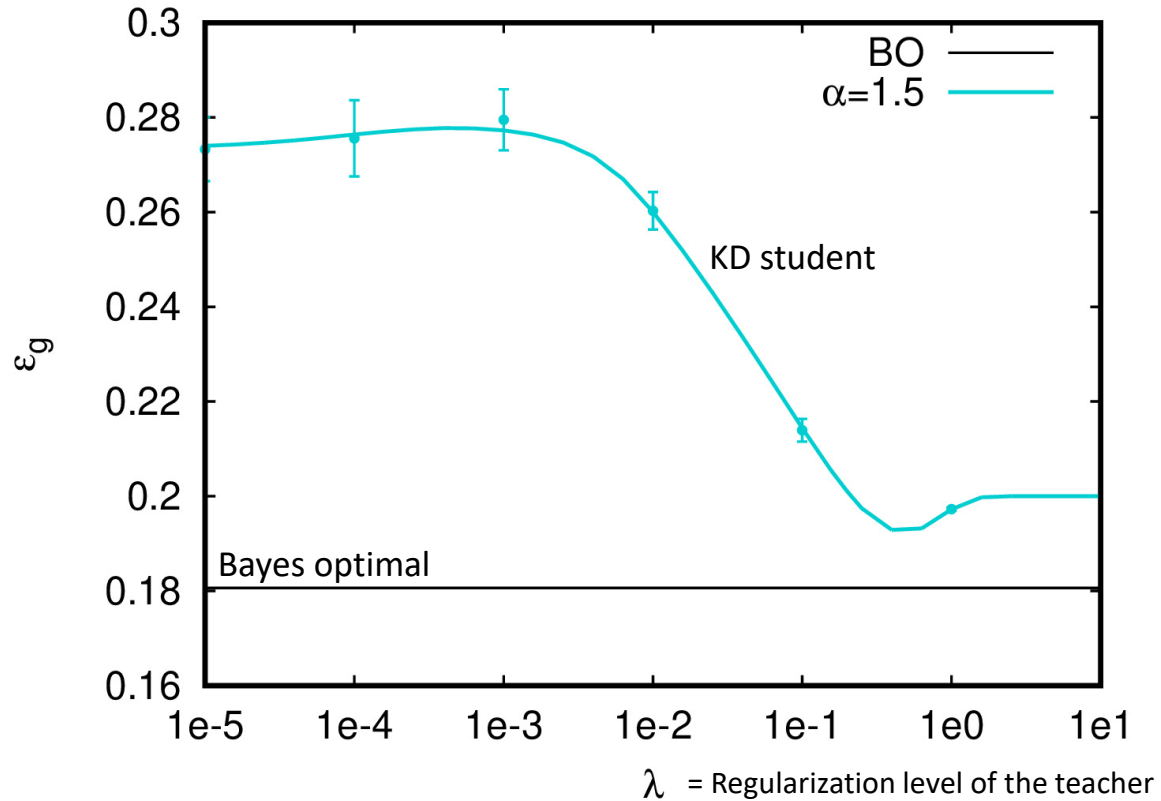
Tuning **regularization intensity λ** is key!!!

Unbalanced clusters: $y^{\mu} = \text{Bern}(\rho)$ **hard!**

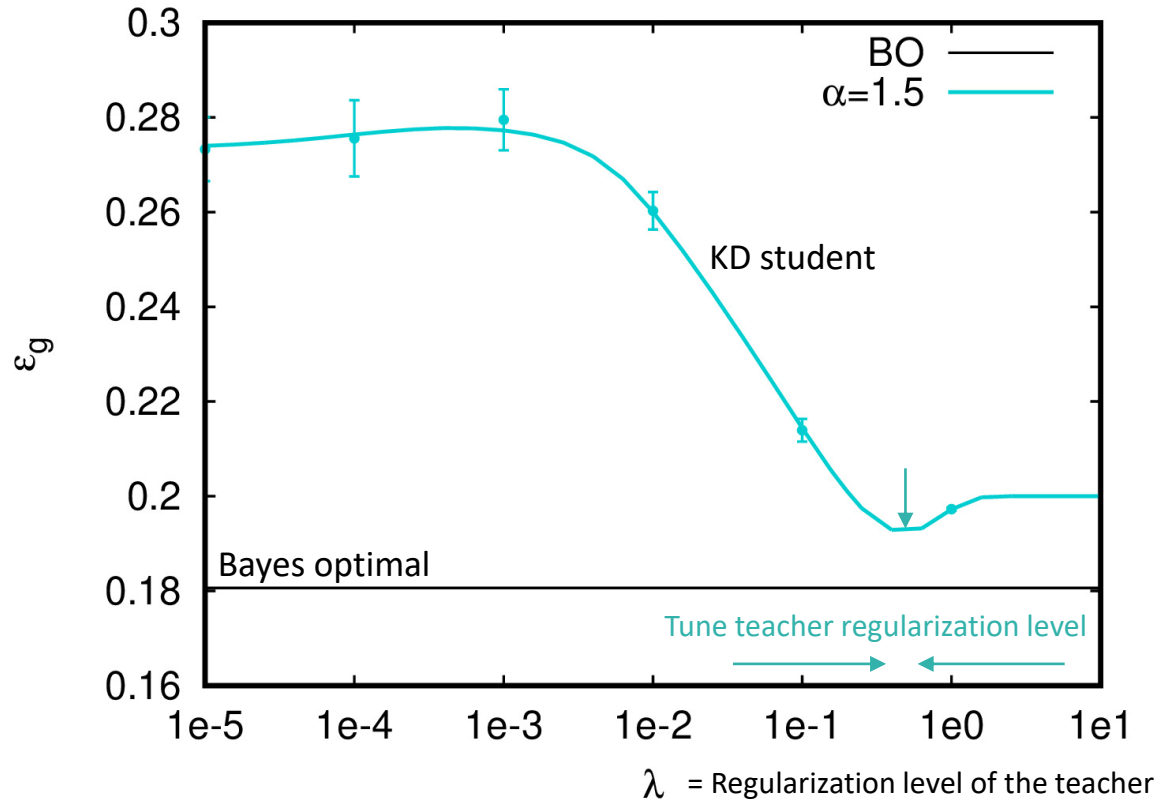


Teacher-student mismatch: weaker student model
Fixed **student sparsity:** fraction $\eta=0.5$ of the weights are trained, the rest set to 0 a priori

L₂-regularized logistic regression teacher:
effect of KD loss on the student



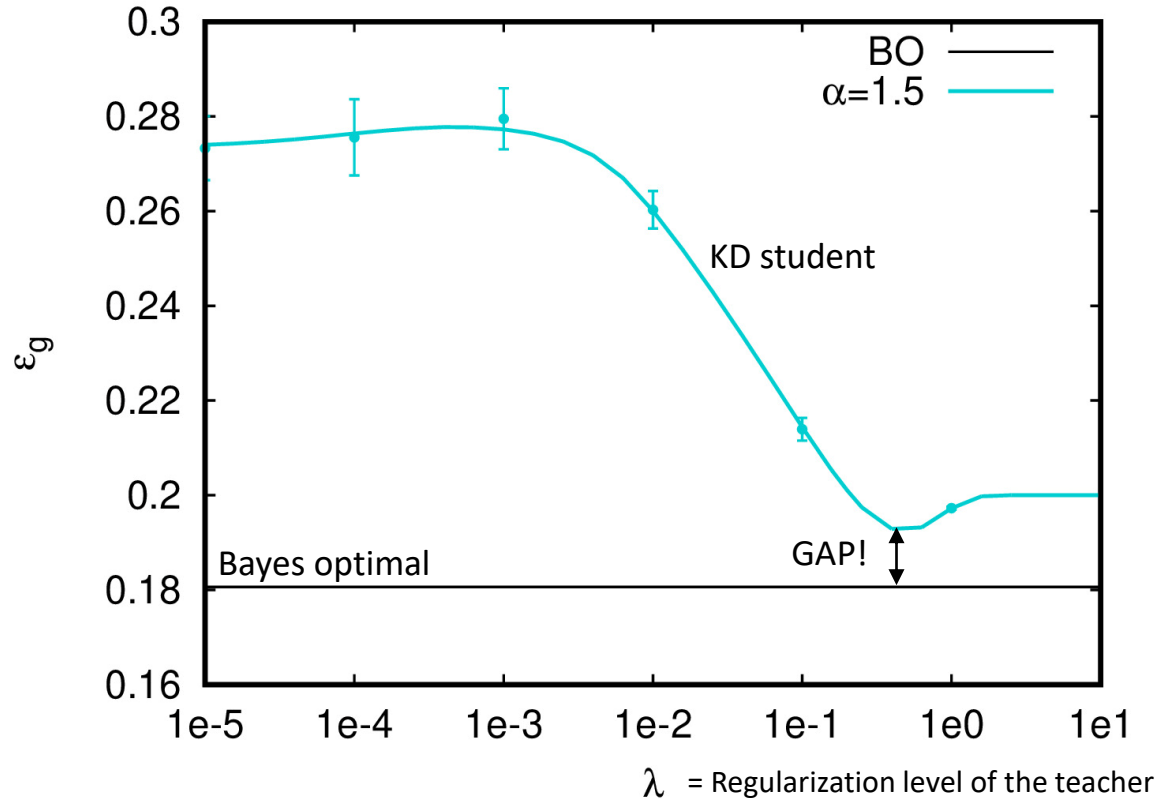
L₂-regularized logistic regression teacher: effect of KD loss on the student



Better teacher=better student

The student inherits the teacher regularization through KD!

L_2 -regularized logistic regression teacher: effect of KD loss on the student



Better teacher=better student

The student inherits the teacher regularization through KD!

The obtained generalization performance is still sub-optimal!

**The KD student improves
together with the teacher**

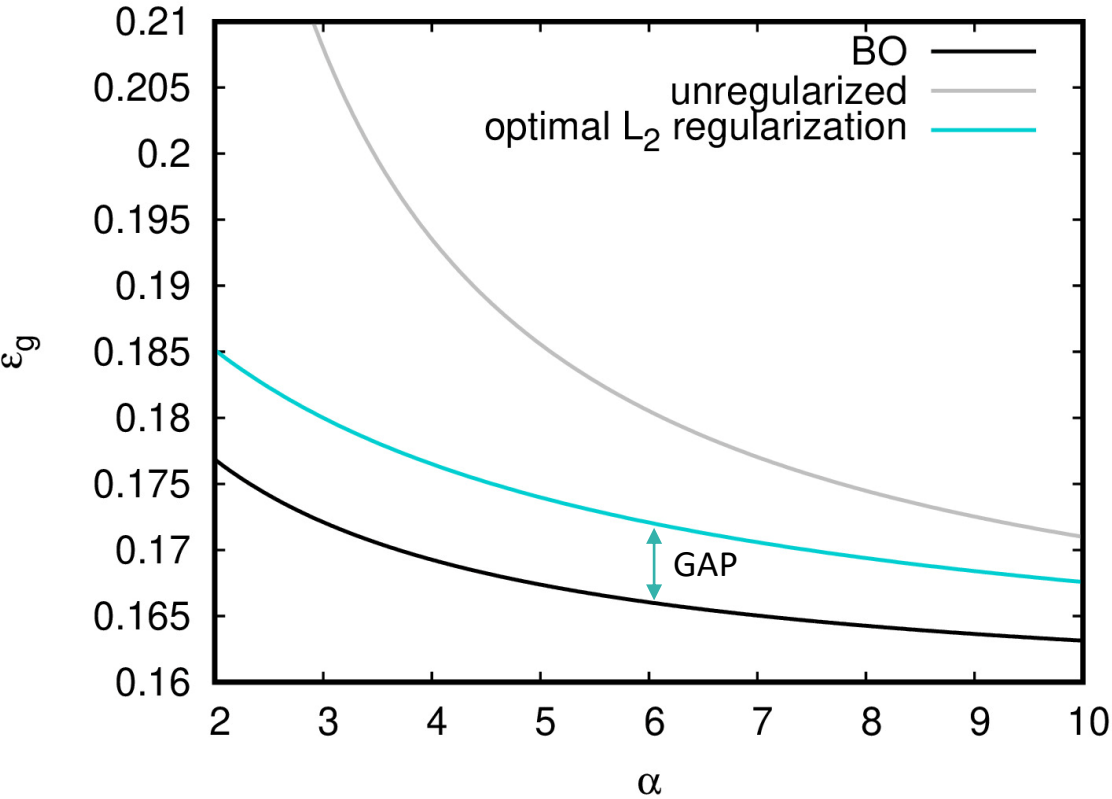
With a **sub-optimal teacher** the
student remains sub-optimal
(as the logistic regression estimator)



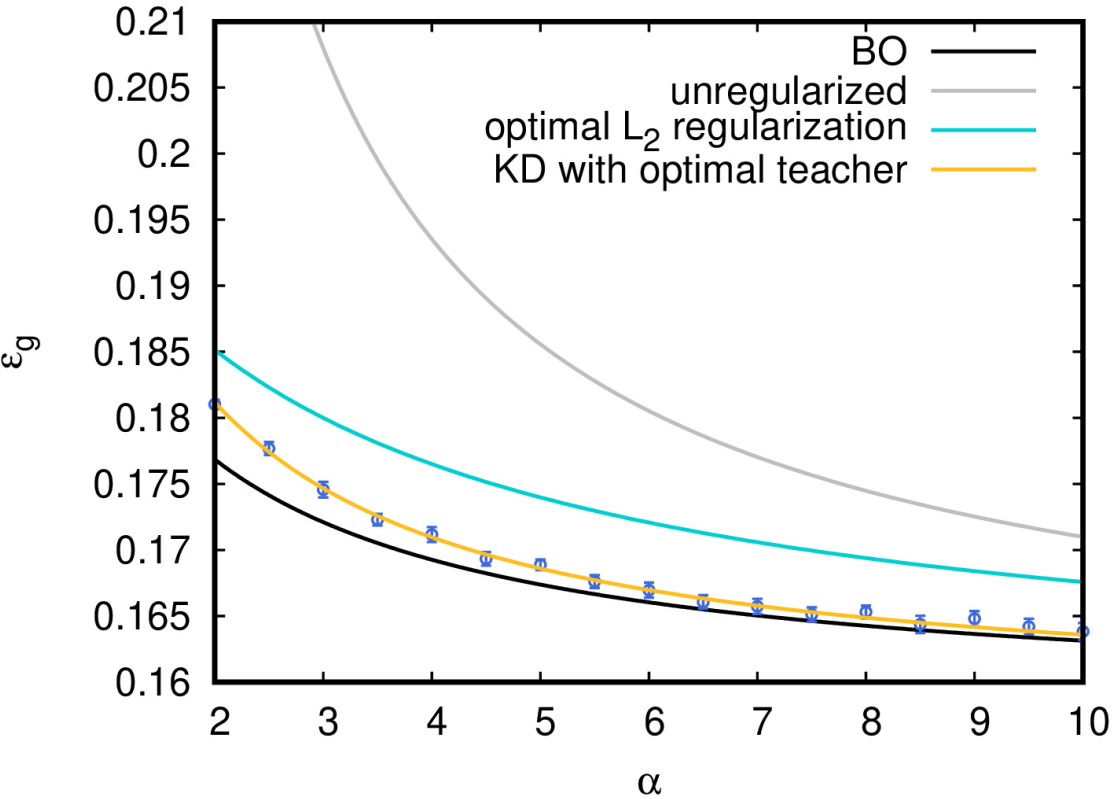
What if the teacher is **not just regularized “explicitly”**?

Is **KD still effective** in transferring
the generalization properties?

Bayes-Optimal teacher:
KD better than logistic regression?



Bayes-Optimal teacher: KD better than logistic regression?



The KD student closes the generalization gap!!!

even though the form of regularization of the teacher is not known explicitly!
DEEP LEARNING READY

TAKE-HOME MESSAGE

With **Knowledge Distillation** the student can inherit the teacher regularization properties:

CONS

Cannot beat an explicit regularization of the same type!

TAKE-HOME MESSAGE

With **Knowledge Distillation** the student can inherit the teacher regularization properties:

CONS

Cannot beat an explicit regularization of the same type!

PROS

Achieve better generalization **even when**
the **form of regularization of the teacher is not known explicitly**
DEEP LEARNING READY :)

TAKE-HOME MESSAGE

With **Knowledge Distillation** the student can inherit the teacher regularization properties:

CONS

Cannot beat an explicit regularization of the same type!

PROS

Achieve better generalization **even when**
the **form of regularization of the teacher is not known explicitly**
DEEP LEARNING READY :)

Thank you for your attention!