Some Recent Insights on Transfer Learning

 $P + Q \rightarrow Q?$

Samory Kpotufe Columbia University, Statistics

Earlier work with G. Martinet, and ongoing work with S. Hanneke

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

Transfer Learning:

Given data $\{X_i, Y_i\} \sim_{i.i.d.} P$, produce a classifier for $(X, Y) \sim Q$.

Case study: Apple Siri's voice assistant

- Initially trained on data from American English speakers ...

- Could not understand 30M+ nonnative speakers in the US!



(日) (四) (日) (日) (日)

Costly Solution \equiv **5+ years acquiring more data and retraining**!

A Main Practical Goal: Cheaply transfer ML software between related populations.

Transfer Learning:

Given data $\{X_i, Y_i\} \sim_{i.i.d.} P$, produce a classifier for $(X, Y) \sim Q$.

Case study: Apple Siri's voice assistant

- Initially trained on data from American English speakers ...

- Could not understand 30M+ nonnative speakers in the US!



A D N A 目 N A E N A E N A B N A C N

Costly Solution \equiv 5+ years acquiring more data and retraining!

A Main Practical Goal:

Cheaply **transfer** ML software between related populations.

Transfer Learning:

Given data $\{X_i, Y_i\} \sim_{i.i.d.} P$, produce a classifier for $(X, Y) \sim Q$.

Case study: Apple Siri's voice assistant

- Initially trained on data from American English speakers ...

- Could not understand 30M+ nonnative speakers in the US!



Costly Solution \equiv 5+ years acquiring more data and retraining!

A Main Practical Goal:

Cheaply transfer ML software between related populations.

AI for Judicial Systems

- Source Population: prison inmates
- Target Population: everyone arrested



イロト 不得 トイヨト イヨト

-

Over 60% inaccurate risk assessments on minorities (2016 Pro-Publica study)

Main Issue: Good Target data is hard or expensive to acquire AI in medicine, Genomics, Insurance Industry, Smart cities,

AI for Judicial Systems

- Source Population: prison inmates
- Target Population: everyone arrested



(日) (四) (日) (日) (日)

Over 60% inaccurate risk assessments on minorities (2016 Pro-Publica study)

Main Issue: Good Target data is hard or expensive to acquire AI in medicine, Genomics, Insurance Industry, Smart cities,

AI for Judicial Systems

- Source Population: prison inmates
- Target Population: everyone arrested



(日) (四) (日) (日) (日)

Over 60% inaccurate risk assessments on minorities (2016 Pro-Publica study)

Main Issue: Good Target data is hard or expensive to acquire AI in medicine, Genomics, Insurance Industry, Smart cities,

AI for Judicial Systems

- Source Population: prison inmates
- Target Population: everyone arrested



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Over 60% inaccurate risk assessments on minorities (2016 Pro-Publica study)

Main Issue: Good Target data is hard or expensive to acquire Al in medicine, Genomics, Insurance Industry, Smart cities,



Many heuristics ... but theory and principles are still evolving

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

What's the relative statistical value of P and Q data?Depends on how far P is from Q ...

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ - つへつ

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

What's the relative statistical value of P and Q data?Depends on how far P is from Q ...

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

What's the relative statistical value of P and Q data?Depends on how far P is from Q ...

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

What's the relative statistical value of P and Q data? Depends on how far P is from Q ...

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

What's the relative statistical value of P and Q data? Depends on how far P is from Q ...

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

What's the relative statistical value of P and Q data?

Depends on how far P is from Q ...

Suppose: \hat{h} is trained on source data $\sim P$, to be *transferred* to target Q.

- Is there sufficient information in source P about target Q?
- If not, how much new data should be collected?
- Would unlabeled data help?
- What's the right mix of P and Q data w.r.t. \$\$ sampling costs?

 $\frac{\text{What's the relative statistical value of } P \text{ and } Q \text{ data?}}{\text{Depends on how } \textit{far } P \text{ is from } Q \dots}$

A D N A 目 N A E N A E N A B N A C N

Formal Setup: Classification $X \mapsto Y$, fixed VC class \mathcal{H}

Given: source data $\{X_i, Y_i\} \sim P^{n_P}$, target data $\{X_i, Y_i\} \sim Q^{n_Q}$. **Goal:** $\hat{h} \in \mathcal{H}$ with small *excess* target error

$$\mathcal{E}_{Q}(\hat{h}) = \mathbb{E}_{Q}\left[\hat{h}(X) \neq Y\right] - \inf_{h \in \mathcal{H}} \mathbb{E}_{Q}\left[h(X) \neq Y\right]$$

Basic Information-theoretic Question: Which $\mathcal{E}_Q(\hat{h})$ is achievable in terms of sample sizes n_P and n_Q ?

Which notion of $dist(P \rightarrow Q)$ captures this error?

Nonparametric work

- (Covariate Shift) [Kpo. and Martinet]
- (Posterior Drift) [Scott 19] [Cai and Wei, 19]
- (Covariate Shift, Posterior Drift) [Reeve, Cannings, Samworth, 21]

Formal Setup:

Classification $X \mapsto Y$, fixed VC class \mathcal{H}

Given: source data $\{X_i, Y_i\} \sim P^{n_P}$, target data $\{X_i, Y_i\} \sim Q^{n_Q}$. Goal: $\hat{h} \in \mathcal{H}$ with small excess target error $\mathcal{E}_Q(\hat{h}) = \mathbb{E}_Q[\hat{h}(X) \neq Y] - \inf \mathbb{E}_Q[h(X) \neq Y]$

Basic Information-theoretic Question: Which $\mathcal{E}_Q(\hat{h})$ is achievable in terms of sample sizes n_P and n_Q ?

Which notion of dist $(P \rightarrow Q)$ captures the drag \mathbb{R}^{2}

Formal Setup: Classification $X \mapsto Y$, fixed VC class \mathcal{H}

Given: source data $\{X_i, Y_i\} \sim P^{n_P}$, target data $\{X_i, Y_i\} \sim Q^{n_Q}$.

Goal: $h \in \mathcal{H}$ with small *excess* target error

$$\mathcal{E}_{Q}(\hat{h}) = \mathbb{E}_{Q}\left[\hat{h}(X) \neq Y\right] - \inf_{h \in \mathcal{H}} \mathbb{E}_{Q}\left[h(X) \neq Y\right]$$

Basic Information-theoretic Question: Which $\mathcal{E}_Q(\hat{h})$ is achievable in terms of sample sizes n_P and n_Q ?

Which notion of $dist(P \rightarrow Q)$ captures this error?

Formal Setup: Classification $X \mapsto Y$, fixed VC class \mathcal{H}

Given: source data $\{X_i, Y_i\} \sim P^{n_P}$, target data $\{X_i, Y_i\} \sim Q^{n_Q}$. Goal: $\hat{h} \in \mathcal{H}$ with small excess target error

$$\mathcal{E}_{Q}(\hat{h}) = \mathbb{E}_{Q}\left[\hat{h}(X) \neq Y\right] - \inf_{h \in \mathcal{H}} \mathbb{E}_{Q}\left[h(X) \neq Y\right]$$

Basic Information-theoretic Question: Which $\mathcal{E}_Q(\hat{h})$ is achievable in terms of sample sizes n_P and n_Q ?

Which notion of $dist(P \rightarrow Q)$ captures this error?

Formal Setup: Classification $X \mapsto Y$, fixed VC class \mathcal{H}

Given: source data $\{X_i, Y_i\} \sim P^{n_P}$, target data $\{X_i, Y_i\} \sim Q^{n_Q}$. Goal: $\hat{h} \in \mathcal{H}$ with small excess target error

$$\mathcal{E}_{Q}(\hat{h}) = \mathbb{E}_{Q}\left[\hat{h}(X) \neq Y\right] - \inf_{h \in \mathcal{H}} \mathbb{E}_{Q}\left[h(X) \neq Y\right]$$

Basic Information-theoretic Question: Which $\mathcal{E}_Q(\hat{h})$ is achievable in terms of sample sizes n_P and n_Q ?

Which notion of $dist(P \rightarrow Q)$ captures this error?

Formal Setup: Classification $X \mapsto Y$, fixed VC class \mathcal{H}

Given: source data $\{X_i, Y_i\} \sim P^{n_P}$, target data $\{X_i, Y_i\} \sim Q^{n_Q}$. Goal: $\hat{h} \in \mathcal{H}$ with small excess target error

$$\mathcal{E}_{Q}(\hat{h}) = \mathbb{E}_{Q}\left[\hat{h}(X) \neq Y\right] - \inf_{h \in \mathcal{H}} \mathbb{E}_{Q}\left[h(X) \neq Y\right]$$

Basic Information-theoretic Question: Which $\mathcal{E}_Q(\hat{h})$ is achievable in terms of sample sizes n_P and n_Q ?

Which notion of $dist(P \rightarrow Q)$ captures this error?

Many competing notions of $dist(P \rightarrow Q)$...

- Extensions of TV: consider |P(A) Q(A)| over suitable A(e.g. d_A divergence/ \mathcal{Y} -discrepancy of S. Ben David, M. Mohri, ...) $\mathcal{E}_Q(\hat{h}) \lesssim o_P(1) + \operatorname{dist}(P \to Q)$
- **Density Ratios:** consider ratio dQ/dP over data space (e.g., Sugiyama, Belkin, Jordan, Wainwright, ...)

 $\mathcal{E}_Q(\hat{h}) \lesssim o_P(1) + \text{estimation error}(d_Q/d_P)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Many competing notions of $dist(P \rightarrow Q)$...



- Extensions of TV: consider |P(A) Q(A)| over suitable A(e.g. d_A divergence/ \mathcal{Y} -discrepancy of S. Ben David, M. Mohri, ...) $\mathcal{E}_Q(\hat{h}) \lesssim o_P(1) + \operatorname{dist}(P \to Q)$
- **Density Ratios:** consider ratio dQ/dP over data space (e.g., Sugiyama, Belkin, Jordan, Wainwright, ...)

 $\mathcal{E}_Q(\hat{h}) \lesssim o_P(1) + \text{estimation error}(d_Q/d_P)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Many competing notions of $dist(P \rightarrow Q)$...



- Extensions of TV: consider |P(A) Q(A)| over suitable A(e.g. d_A divergence/ \mathcal{Y} -discrepancy of S. Ben David, M. Mohri, ...) $\mathcal{E}_Q(\hat{h}) \lesssim o_P(1) + \operatorname{dist}(P \to Q)$
- **Density Ratios:** consider ratio dQ/dP over data space (e.g., Sugiyama, Belkin, Jordan, Wainwright, ...)

 $\mathcal{E}_Q(\hat{h}) \lesssim o_P(1) + \text{estimation error}(d_Q/d_P)$

Namely: P far from $Q \implies$ Transfer is Hard



Namely: P far from $Q \implies$ Transfer is Hard

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

Namely: P far from $Q \implies$ Transfer is Hard

Namely: P far from $Q \implies$ Transfer is Hard



Large TV, d_A , \mathcal{Y} -disc $\approx 1/2$

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Namely: P far from $Q \implies$ Transfer is Hard



Asymmetry in transfer \implies Metrics are inappropriate

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Namely: P far from $Q \implies$ Transfer is Hard

Source Distribution

Target Distribution

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00



Large dQ/dP, KL-div $\approx \infty$

Intuition: $h \in \mathcal{H}$ has low error under $P \implies$ low error under QFor now assume $h_P^* = h_Q^* \dots$ Transfer exponent $\rho > 0$: $\forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h) \leq c \cdot \mathcal{E}_P^{1/\rho}(h)$

ho captures a continuum of easy to hard transfer

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Intuition: $h \in \mathcal{H}$ has low error under $P \implies$ low error under Q

For now assume $h_P^* = h_Q^* \dots$

Transfer exponent $\rho > 0$: $\forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h) \leq c \cdot \mathcal{E}_P^{1/\rho}(h)$

ho captures a continuum of easy to hard transfer

Intuition: $h \in \mathcal{H}$ has low error under $P \implies$ low error under QFor now assume $h_P^* = h_Q^* \dots$ Transfer exponent $\rho > 0$: $\forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h) \leq c \cdot \mathcal{E}_P^{1/\rho}(h)$

ho captures a continuum of easy to hard transfer

Intuition: $h \in \mathcal{H}$ has low error under $P \implies$ low error under Q

For now assume $h_P^* = h_Q^* \dots$

 $\begin{array}{l} \text{Transfer exponent } \rho > 0 \text{:} \\ \forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h) \leq c \cdot \mathcal{E}_P^{1/\rho}(h) \end{array}$

ho captures a continuum of easy to hard transfer ...

Intuition: $h \in \mathcal{H}$ has low error under $P \implies$ low error under Q

For now assume $h_P^* = h_Q^* \ldots$

 $\begin{array}{l} \text{Transfer exponent } \rho > 0 \text{:} \\ \forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h) \leq c \cdot \mathcal{E}_P^{1/\rho}(h) \end{array}$

 ρ captures a continuum of easy to hard transfer ...

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
$extsf{Transfer exponent } ho > 0 extsf{:}$ $orall h \in \mathcal{H}, \quad \mathcal{E}_Q(h,h^*) \leq c \cdot \mathcal{E}_P^{1/ ho}(h,h^*)$

 $\begin{array}{l} \text{Transfer exponent } \rho > 0 \text{:} \\ \forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h,h^*) \leq c \cdot \mathcal{E}_P^{1/\rho}(h,h^*) \end{array}$

For deterministic $Y = h^*(X)$ this reduces to:

 $Q_X(h \neq h^*) \le c \cdot P_X^{1/\rho}(h \neq h^*)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Transfer exponent $\rho > 0$:

 $\forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h, h^*) \le c \cdot \mathcal{E}_P^{1/\rho}(h, h^*)$



 $\rho = 1$ but $d_{\mathcal{A}}(P,Q) = \mathcal{Y}\text{-disc}(P,Q) = 1/4$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

$\begin{array}{l} \text{Transfer exponent } \rho > 0 \text{:} \\ \forall h \in \mathcal{H}, \quad \mathcal{E}_Q(h,h^*) \leq c \cdot \mathcal{E}_P^{1/\rho}(h,h^*) \end{array}$



ho=1 but KL, Renyi, blow up ...

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

 $extsf{Transfer exponent }
ho > 0 extsf{:}$ $orall h \in \mathcal{H}, \quad \mathcal{E}_Q(h,h^*) \leq c \cdot \mathcal{E}_P^{1/
ho}(h,h^*)$



 $\rho > 1 \equiv$ how much P covers decision boundary

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

$ext{Transfer exponent } ho > 0 ext{:}$ $orall h \in \mathcal{H}, \quad \mathcal{E}_Q(h,h^*) \leq c \cdot \mathcal{E}_P^{1/ ho}(h,h^*)$



 $0 < \rho < 1 \equiv$ Super Transfer (P has better coverage of decision boundary)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

 ρ captures performance limits (minimax rates) under transfer \ldots

Easy to hard classification





Easy ClassificationHard ClassificationEssential: Noise in Y|X, and X-mass near decision boundaryBernstein condition: $Q_X(h \neq h^*) \lesssim \mathcal{E}_Q{}^\beta(h;h^*), \quad \beta \in [0,1]$

Easy to hard classification





Easy Classification

Hard Classification

Essential: Noise in Y|X, and X-mass near decision boundary **Bernstein condition:** $Q_X(h \neq h^*) \lesssim \mathcal{E}_Q^{\beta}(h;h^*), \quad \beta \in [0,1]$

Easy to hard classification





Easy Classification

Hard Classification

Essential: Noise in Y|X, and X-mass near decision boundary

Bernstein condition: $Q_X(h \neq h^*) \lesssim \mathcal{E}_Q^{\beta}(h;h^*), \quad \beta \in [0,1]$

Easy to hard classification





Easy Classification

Hard Classification

Essential: Noise in Y|X, and X-mass near decision boundary Bernstein condition: $Q_X(h \neq h^*) \lesssim \mathcal{E}_O^{\beta}(h;h^*), \quad \beta \in [0,1]$

Easy to hard classification





Easy Classification

Hard Classification

Essential: Noise in Y|X, and X-mass near decision boundary Bernstein condition: $Q_X(h \neq h^*) \lesssim \mathcal{E}_O^\beta(h;h^*), \quad \beta \in [0,1]$

Given: labeled source and target data $\{X_i, Y_i\} \sim P^{n_P} \times Q^{n_Q}$.

Theorem. Let \hat{h} trained on samples from P + Q:

$$\inf_{\hat{h}} \sup_{(P,Q)} \mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/\rho} + n_Q \right)^{-1/(2-\beta)}$$

Tight for any $\mathcal{H}, \rho \geq 1, \beta, n_P, n_Q \dots$

- Benefits of Unlabeled data: cannot improve the rates ...
- Benefits of Labeled Q data: transition at $n_Q > n_P^{1/
 ho}$
- Adaptive sampling at optimal \$\$ costs: possible in some regimes

Given: labeled source and target data $\{X_i, Y_i\} \sim P^{n_P} \times Q^{n_Q}$.

Theorem. Let h trained on samples from P + Q:

$$\inf_{\hat{h}} \sup_{(P,Q)} \mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/\rho} + n_Q \right)^{-1/(2-\beta)}$$

Tight for any $\mathcal{H}, \rho \geq 1, \beta, n_P, n_Q \dots$

- Benefits of Unlabeled data: cannot improve the rates ...
- Benefits of Labeled Q data: transition at $n_Q > n_P^{1/\rho}$
- Adaptive sampling at optimal \$\$ costs: possible in some regimes

Given: labeled source and target data $\{X_i, Y_i\} \sim P^{n_P} \times Q^{n_Q}$.

Theorem. Let \hat{h} trained on samples from P + Q:

$$\inf_{\hat{h}} \sup_{(P,Q)} \mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/\rho} + n_Q\right)^{-1/(2-\beta)}$$

Tight for any $\mathcal{H}, \rho \geq 1, \beta, n_P, n_Q \dots$

- Benefits of Unlabeled data: cannot improve the rates ...
- Benefits of Labeled Q data: transition at $n_Q > n_P^{1/
 ho}$
- Adaptive sampling at optimal \$\$ costs: possible in some regimes

Given: labeled source and target data $\{X_i, Y_i\} \sim P^{n_P} \times Q^{n_Q}$.

Theorem. Let \hat{h} trained on samples from P + Q:

$$\inf_{\hat{h}} \sup_{(P,Q)} \mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/\rho} + n_Q\right)^{-1/(2-\beta)}$$

Tight for any $\mathcal{H}, \rho \geq 1, \beta, n_P, n_Q$...

- Benefits of Unlabeled data: cannot improve the rates ...
- Benefits of Labeled Q data: transition at $n_Q > n_P^{1/
 ho}$
- Adaptive sampling at optimal \$\$ costs: possible in some regimes

Given: labeled source and target data $\{X_i, Y_i\} \sim P^{n_P} \times Q^{n_Q}$.

Theorem. Let \hat{h} trained on samples from P + Q:

$$\inf_{\hat{h}} \sup_{(P,Q)} \mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/\rho} + n_Q\right)^{-1/(2-\beta)}$$

Tight for any $\mathcal{H}, \rho \geq 1, \beta, n_P, n_Q \dots$

- Benefits of Unlabeled data: cannot improve the rates ...
- Benefits of Labeled Q data: transition at $n_Q > n_P^{1/\rho}$
- Adaptive sampling at optimal \$\$ costs: possible in some regimes

Lower-Bound Analysis

\hat{h} has access to (P,Q) samples, but has to do well on just Q ...

Construction: family $\{(P,Q)_h\}$, any \mathcal{H} , $\rho \geq 1$, β :

- $(P^{n_P} \times Q^{n_Q})_h$ are close in KL-divergence
- But far under distance $Q_h(h' \neq h)$

The rest is extensions of Fano (see e.g. Tsybakov, or Barron and Li) ...

Lower-Bound Analysis

 \hat{h} has access to (P,Q) samples, but has to do well on just Q ...

Construction: family $\{(P,Q)_h\}$, any \mathcal{H} , $\rho \geq 1$, β :

- $(P^{n_P} \times Q^{n_Q})_h$ are close in KL-divergence
- But far under distance $Q_h(h' \neq h)$

The rest is extensions of Fano (see e.g. Tsybakov, or Barron and Li) ...

Lower-Bound Analysis

 \hat{h} has access to (P,Q) samples, but has to do well on just Q ...

Construction: family $\{(P,Q)_h\}$, any \mathcal{H} , $\rho \geq 1$, β :

- $(P^{n_P} \times Q^{n_Q})_h$ are close in KL-divergence
- But far under distance $Q_h(h' \neq h)$

The rest is extensions of Fano (see e.g. Tsybakov, or Barron and Li) ...

Performance limits:
$$\mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/
ho} + n_Q
ight)^{-1/(2-eta)}$$

(Optimal Heuristics for unknown ρ)

Low Classification noise ($\beta = 1$):

ERM on combined source and target data.

Non i.i.d. Bernstein + usual fixed point argument

Unknown Noise Level ($\beta \in [0,1]$):

Minimize $\hat{R}_Q(h)$ subject to $\hat{R}_P(h) \leq \min_{h'} \hat{R}_P(h') + \Delta_{n_P}(h)$

Lepski-type argument

人口 医水黄 医水黄 医水黄素 化甘油

Performance limits:
$$\mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/\rho} + n_Q\right)^{-1/(2-\beta)}$$

We are interested in $\textit{adaptivity} \text{ to } \rho \ldots$

(Optimal Heuristics for unknown ρ)

Low Classification noise ($\beta = 1$):

ERM on combined source and target data.

Non i.i.d. Bernstein + usual fixed point argument

Unknown Noise Level ($\beta \in [0,1]$): Minimize $\hat{R}_Q(h)$ subject to $\hat{R}_P(h) \le \min_{h'} \hat{R}_P(h') + \Delta_{n_P}(h)$

Lepski-type argument

Performance limits:
$$\mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/
ho} + n_Q
ight)^{-1/(2-eta)}$$

(Optimal Heuristics for unknown ρ)

Low Classification noise ($\beta = 1$):

ERM on combined source and target data.

Non i.i.d. Bernstein + usual fixed point argument

Unknown Noise Level ($\beta \in [0,1]$):

Minimize $\hat{R}_Q(h)$ subject to $\hat{R}_P(h) \leq \min_{h'} \hat{R}_P(h') + \Delta_{n_P}(h)$

Lepski-type argument

・ロット (雪) ・ (日) ・ (日) ・ (日)

Performance limits:
$$\mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/
ho} + n_Q\right)^{-1/(2-eta)}$$

(Optimal Heuristics for unknown ρ)

Low Classification noise ($\beta = 1$):

ERM on combined source and target data.

Non i.i.d. Bernstein + usual fixed point argument

Unknown Noise Level ($\beta \in [0, 1]$):

Minimize $\hat{R}_Q(h)$ subject to $\hat{R}_P(h) \leq \min_{h'} \hat{R}_P(h') + \Delta_{n_P}(h)$

Lepski-type argument

A D N A 目 N A E N A E N A B N A C N

Performance limits:
$$\mathcal{E}_Q(\hat{h}) \propto \left(n_P^{1/
ho} + n_Q
ight)^{-1/(2-eta)}$$

(Optimal Heuristics for unknown ρ)

Low Classification noise ($\beta = 1$):

ERM on combined source and target data.

Non i.i.d. Bernstein + usual fixed point argument

Unknown Noise Level ($\beta \in [0, 1]$): Minimize $\hat{R}_Q(h)$ subject to $\hat{R}_P(h) \leq \min_{h'} \hat{R}_P(h') + \Delta_{n_P}(h)$

Lepski-type argument

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

Quick Summary and some New Directions ...

(ロ)、(型)、(E)、(E)、 E) の(()

• ρ captures a more optimistic view of transferability $P \to Q$. • Reveals general form of optimal heuristics:

Minimize $\hat{R}_P(h)$ subject to $\hat{R}_Q(h)$ not too large ...

- Cost-sensitive sampling is possible with no knowledge of ρ .
- Results extend to $h_P^* \neq h_Q^*$: $\exists \hat{h} \text{ s.t.}$

$$\mathcal{E}_Q(\hat{h}) \lesssim \min\left\{n_P^{-1/(2-\beta)\rho} + \mathcal{E}_Q(h_P^*), n_Q^{-1/(2-\beta)}\right\}$$

- ρ captures a more optimistic view of transferability P → Q.
 Reveals general form of optimal heuristics:
 Minimize R̂_P(h) subject to R̂_O(h) not too large ...
- Cost-sensitive sampling is possible with no knowledge of ρ .
- Results extend to $h_P^* \neq h_Q^*$: $\exists \hat{h} \text{ s.t.}$

$$\mathcal{E}_Q(\hat{h}) \lesssim \min\left\{n_P^{-1/(2-\beta)\rho} + \mathcal{E}_Q(h_P^*), n_Q^{-1/(2-\beta)}\right\}$$

- ρ captures a more optimistic view of transferability $P \rightarrow Q$.
- Reveals general form of optimal heuristics:

Minimize $\hat{R}_P(h)$ subject to $\hat{R}_Q(h)$ not too large ...

- Cost-sensitive sampling is possible with no knowledge of ρ .
- Results extend to $h_P^* \neq h_Q^*$: $\exists \hat{h} \text{ s.t.}$

$$\mathcal{E}_Q(\hat{h}) \lesssim \min\left\{n_P^{-1/(2-\beta)\rho} + \mathcal{E}_Q(h_P^*), n_Q^{-1/(2-\beta)}\right\}$$

- ρ captures a more optimistic view of transferability $P \rightarrow Q$.
- Reveals general form of optimal heuristics:

Minimize $\hat{R}_P(h)$ subject to $\hat{R}_Q(h)$ not too large ...

Cost-sensitive sampling is possible with no knowledge of ρ.

• Results extend to $h_P^* \neq h_Q^*$: $\exists \hat{h} \text{ s.t.}$

 $\mathcal{E}_Q(\hat{h}) \lesssim \min\left\{n_P^{-1/(2-\beta)\rho} + \mathcal{E}_Q(h_P^*), n_Q^{-1/(2-\beta)}\right\}$

- ρ captures a more optimistic view of transferability $P \rightarrow Q$.
- Reveals general form of optimal heuristics:

Minimize $\hat{R}_P(h)$ subject to $\hat{R}_Q(h)$ not too large ...

- Cost-sensitive sampling is possible with no knowledge of ρ .
- Results extend to $h_P^* \neq h_Q^*$: $\exists \hat{h} \text{ s.t.}$

$$\mathcal{E}_Q(\hat{h}) \lesssim \min\left\{n_P^{-1/(2-\beta)\rho} + \mathcal{E}_Q(h_P^*), n_Q^{-1/(2-\beta)}\right\}$$

Current work:

Performance limits in multi-task (new arXiv with S. Hanneke)

$P_1 + P_2 + \dots + P_N + Q \to Q?$

Prior theory only yields single source rates ...

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

Current work:

Performance limits in multi-task (new arXiv with S. Hanneke)

$P_1 + P_2 + \dots + P_N + Q \to Q?$

Prior theory only yields single source rates ...

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Setup:

$N \text{ sources } \{P_t\}_{t=1}^N \mapsto Q \text{ with } \mathcal{E}_t(h) \gtrsim \mathcal{E}_Q^{\rho_t}(h)$

Minimax Rate on
$$\mathcal{E}_Q(\hat{h})$$
 : $\min_{t\in [N+1]} \left(\sum_{s=1}^t n_{(t)}\right)^{-1/(2-eta)ar{
ho}_t}$

Adaptive Strategies (as $N \to \infty$):

Low noise ($\beta = 1$): ERM on combined data

Information on ranking $ho_{(1)} \leq ... \leq
ho_{(N)}$: Greedy strategy ...

No adaptive strategy outside above regimes !!!

$$\inf_{\hat{h}} \sup_{\text{rankings } \{P_t\} \times Q} \mathcal{E}_Q(\hat{h}) \gtrsim n_Q^{-1/(2-\beta)}$$

Setup:

N sources $\{P_t\}_{t=1}^N \mapsto Q$ with $\mathcal{E}_t(h) \gtrsim \mathcal{E}_Q^{\rho_t}(h)$

Minimax Rate on
$$\mathcal{E}_Q(\hat{h})$$
 : $\min_{t\in [N+1]} \left(\sum_{s=1}^t n_{(t)}\right)^{-1/(2-eta)ar{
ho}_t}$

Adaptive Strategies (as $N \to \infty$):

Low noise $(\beta = 1)$: ERM on combined data

Information on ranking $ho_{(1)} \leq ... \leq
ho_{(N)}$: Greedy strategy ...

No adaptive strategy outside above regimes !!!

$$\inf_{\hat{h}} \sup_{\text{rankings } \{P_t\} \times Q} \mathcal{E}_Q(\hat{h}) \gtrsim n_Q^{-1/(2-\beta)}$$

Setup:

N sources $\{P_t\}_{t=1}^N \mapsto Q$ with $\mathcal{E}_t(h) \gtrsim \mathcal{E}_Q^{\rho_t}(h)$

$$\text{Minimax Rate on } \mathcal{E}_Q(\hat{h}) : \quad \min_{t \in [N+1]} \left(\sum_{s=1}^t n_{(t)} \right)^{-1/(2-\beta)\bar{\rho}_t}$$

Adaptive Strategies (as $N \to \infty$):

Low noise $(\beta = 1)$: ERM on combined data

Information on ranking $ho_{(1)} \leq ... \leq
ho_{(N)}$: Greedy strategy ...

No adaptive strategy outside above regimes !!!

$$\inf_{\hat{h}} \sup_{\text{rankings } \{P_t\} \times Q} \mathcal{E}_Q(\hat{h}) \gtrsim n_Q^{-1/(2-\beta)}$$
Setup:

N sources $\{P_t\}_{t=1}^N \mapsto Q$ with $\mathcal{E}_t(h) \gtrsim \mathcal{E}_Q^{\rho_t}(h)$

$$\mathsf{Minimax} \; \mathsf{Rate} \; \mathsf{on} \; \mathcal{E}_Q(\hat{h}) : \quad \min_{t \in [N+1]} \left(\sum_{s=1}^t n_{(t)} \right)^{-1/(2-\beta)\bar{\rho}_t}$$

Adaptive Strategies (as $N \to \infty$):

Low noise ($\beta = 1$): ERM on combined data

Information on ranking $ho_{(1)} \leq ... \leq
ho_{(N)}$: Greedy strategy ...

No adaptive strategy outside above regimes !!!

$$\inf_{\hat{h}} \sup_{\text{rankings } \{P_t\} \times Q} \sup_{Q} \mathcal{E}_Q(\hat{h}) \gtrsim n_Q^{-1/(2-\beta)}$$

Setup:

N sources $\{P_t\}_{t=1}^N \mapsto Q$ with $\mathcal{E}_t(h) \gtrsim \mathcal{E}_Q^{\rho_t}(h)$

$$\mathsf{Minimax} \; \mathsf{Rate} \; \mathsf{on} \; \mathcal{E}_Q(\hat{h}) : \quad \min_{t \in [N+1]} \left(\sum_{s=1}^t n_{(t)} \right)^{-1/(2-\beta)\bar{\rho}_t}$$

Adaptive Strategies (as $N \to \infty$):

Low noise ($\beta = 1$): ERM on combined data

Information on ranking $\rho_{(1)} \leq ... \leq \rho_{(N)}$: Greedy strategy ...

No adaptive strategy outside above regimes !!!

$$\inf_{\hat{h}} \sup_{\text{rankings } \{P_t\} \times Q} \sup_{Q} \mathcal{E}_Q(\hat{h}) \gtrsim n_Q^{-1/(2-\beta)}$$

Setup:

 $N \text{ sources } \{P_t\}_{t=1}^N \mapsto Q \text{ with } \mathcal{E}_t(h) \gtrsim \mathcal{E}_Q^{\rho_t}(h)$

$$\mathsf{Minimax} \; \mathsf{Rate} \; \mathsf{on} \; \mathcal{E}_Q(\hat{h}) : \qquad \min_{t \in [N+1]} \left(\sum_{s=1}^t n_{(t)} \right)^{-1/(2-\beta)\bar{\rho}_t}$$

Adaptive Strategies (as $N \to \infty$):

Low noise ($\beta = 1$): ERM on combined data

Information on ranking $\rho_{(1)} \leq ... \leq \rho_{(N)}$: Greedy strategy ...

No adaptive strategy outside above regimes !!!

$$\inf_{\hat{h}} \sup_{\text{rankings}} \sup_{\{P_t\} \times Q} \mathcal{E}_Q(\hat{h}) \gtrsim n_Q^{-1/(2-\beta)}.$$

Somehow we are still just scratching the surface ...

Thanks!

Somehow we are still just scratching the surface ...

Thanks!