



Figure 4: Normalized L_2 distance between two bivariate distributions with standard normal marginals, one with Gaussian copula and the other with maximum-entropy copula under Spearman correlation constraints, for various Spearman correlation values ρ_s .

A Further Details

A.1 Additional Experiment

It can be shown that the highest R^2 and the lowest RMSE achievable by a regression model using \mathbf{x} to predict y read respectively

$$\bar{R}^2(P_{\mathbf{x},y}) = 1 - e^{-2I(y;\mathbf{x})}$$

and

$$RMSE(P_{\mathbf{x},y}) = e^{-I(y;\mathbf{x})} \sqrt{\text{Var}(y)}.$$

We use MIND to estimate the highest performance achievable in the Kaggle competition ‘House Prices: Advanced Regression Techniques’. The aim is to predict the price at which houses were sold from various continuous and categorical variables. The results are summarized in Table 3, when all variables are used, when the 10 explanatory variables OverallQual, GrLivArea, YearBuilt, TotalBsmtSF, OverallCond, LotArea, BsmtFinSF1, BldgType, KitchenQual, MSZoning are used, and when the first 5 of the foregoing list are used.

A.2 Further Details on Handling Categorical Data

Categorical and non-ordinal variables should be ordinarily encoded as customary, and ordinal data should be treated as continuous variables. The only practical requirements for the validity of this approach are i) to use a ranking function that assigns different ranks to all inputs including ties (e.g. `scipy`’s ‘`rankdata`’ function with method ‘`ordinal`’, or leveraging `PyTorch`’s or `Tensorflow`’s ‘`argsort`’), and ii) to avoid encoding

methods that may result in linearly dependent coordinates (e.g. one-hot-encoding on a binary non-ordinal categorical variable). When a suitable ranking function is not available a small random jitter may be added to ordinal variables to remove ties.

This approach is mathematically valid thanks to the quantization characterization of the mutual information (Cover (1999), Definition 8.54). In effect, if we denote \mathcal{P} (resp. \mathcal{Q}) a partition of the domain \mathcal{X} (resp. \mathcal{Y}) of \mathbf{x} (resp. \mathbf{y}), and $[\mathbf{x}]_{\mathcal{P}} \in \mathbb{N}$ (resp. $[\mathbf{y}]_{\mathcal{Q}} \in \mathbb{N}$) a discrete random variable indicating which element of \mathcal{P} (resp. \mathcal{Q}) \mathbf{x} (resp. \mathbf{y}) belongs to, then we have

$$I(\mathbf{y}; \mathbf{x}) = \sup_{\mathcal{P}, \mathcal{Q}} I([\mathbf{y}]_{\mathcal{Q}}; [\mathbf{x}]_{\mathcal{P}}), \quad (12)$$

where the rightmost mutual information is between discrete random variables.

This characterization implies that the mutual information is always invariant by 1-to-1 transformations, of which ordinal encoding is one, whether coordinates are all continuous, all categorical, or a mix.⁵ It can also be seen that adding a negligible random jitter to an ordinal random variable will not materially change the mutual information,⁶ but will turn the ordinal variable into a continuous one so that results developed for continuous variables may apply. Strictly speaking, by the data processing inequality (Cover (1999), Theorem 2.8.1), adding a random jitter to ordinal variables increases the mutual information but, as the jitter standard deviation goes to zero, the difference becomes negligible, even though we still enjoy the benefits of working with continuous variables, without downside. Considering that MIND only depends on variables through their ranks, we may do without adding a jitter, so long as the ranking algorithm does not attribute the same rank to ties.

Another approach for handling categorical variables would be to use Table 2, and to use the entropy decomposition formula (Equation (1)) to estimate differential entropies and conditional entropies. However, this approach can be far less data-efficient as it requires splitting the dataset into as many subsets as the number of distinct tuples of categorical variable values, so as to evaluate conditional differential entropies.

⁵Indeed, for any 1-to-1 transformation f and partition \mathcal{P} of \mathcal{X} , we may always find a partition \mathcal{P}' of the image space $f(\mathcal{X})$ such that $[\mathbf{x}]_{\mathcal{P}} = [f(\mathbf{x})]_{\mathcal{P}'}$.

⁶If we denote g the operation consisting of adding a negligible random jitter to an ordinal variable x_i , then by reducing the jitter’s standard deviation, for any partition \mathcal{P} of the domain of x_i , we may always find a partition \mathcal{P}' of the image space such that $\mathbb{P}([x_i]_{\mathcal{P}} = [g(x_i)]_{\mathcal{P}'}) = 1 - e$ for e arbitrarily small.

	y is continuous	y is categorical
x is continuous	$I(y; x) = h(y) + h(x) - h(y; x)$	$I(y; x) = h(x) - \sum_{i \in \mathcal{Y}} h(x y=i) P_y(i)$
x is categorical	$I(y; x) = h(y) - \sum_{i \in \mathcal{X}} h(y x=i) P_x(i)$	$I(y; x) = H(y) + H(x) - H(y; x)$
x has continous coordinates x_c and categorical coordinates x_d	$I(y; x) = h(y) + \sum_{i \in \mathcal{X}_d} [h(x_c x_d=i) - h(y, x_c x_d=i)] P_{x_d}(i)$	$I(y; x) = I(y; x_d) + \sum_{i \in \mathcal{X}_d} P_{x_d}(i) h(x_c x_d=i) - \sum_{j \in \mathcal{Y}} h(x_c x_d=i, y=j) P_{x_d, y}(i, j)$

Table 2: Expression of the mutual information $I(y; x)$ as a function of the Shannon entropy $H(\cdot)$, and/or the differential entropy $h(\cdot)$, depending on whether y and/or x has continuous and/or categorical coordinates. Expressions of the type $h(x|y=i)$ are to be understood as the differential entropy of the continuous conditional distribution $x|y=i$.

$I(y; x)$	\bar{R}^2	$RMSE$	d	n
1.50	0.95	18,531	80	1460
0.76	0.78	36,979	10	1460
0.65	0.73	41,007	5	1460

Table 3: Mutual information and highest performances achievable in the ‘House Prices: Advanced Regression Techniques’ Kaggle challenge using all, 10 or 5 explanatory variables.

B Proofs

B.1 Proof of Proposition 2.1

Any continuous 1-to-1 univariate transformation g_i is either increasing or decreasing. Increasing transformations leave copulas invariant. Moreover, any continuous decreasing function g_i on \mathbb{R} can be written as $f(-x)$ where f is a continuous increasing function, so that we may focus on proving the statement for the transformation $g_i : z \rightarrow -z$. Let us denote $\bar{z} = (z_1, \dots, -z_i, \dots, z_d)$, \bar{c} its copula density, and c the copula density of $z = (z_1, \dots, z_d)$. We have

$$\bar{c}(u_1, \dots, u_d) = c(u_1, \dots, 1 - u_i, \dots, u_d).$$

A simple change of variables shows that $h(u_z) = h(u_{\bar{z}})$.

B.2 Proof of Theorem 3.1

Let P be a d -dimensional copula distribution, and U the uniform distribution on $[0, 1]^d$. We note that $h(P) = -KL(P||U)$. Thus the optimization problem (MIND) is equivalent to looking for the I -projection of U on the space \mathcal{E} of copula distributions satisfying the linear constraint $E_P[\phi_m(u)] = \alpha_m$, as defined in Csiszár (1975).

Existence and uniqueness: If there exists a copula distribution P satisfying the constraints and admitting an entropy, then \mathcal{E} is not empty. \mathcal{E} is convex as every convex combination of copulas satisfying the linear constraint $E_P[\phi_m(u)] = \alpha_m$ is itself a copula that satisfies said constraint.

We say that a space of continuous distributions supported on $[0, 1]^d$ is variation closed when it is closed in the topology of the variation distance $|P - Q| = \int |p - q| dU$, where p and q are the Radon-Nikodym derivatives of P and Q with respect to U (i.e. their pdfs).

Lemma B.1. \mathcal{E} is variation-closed.

Proof. Let $P_n \in \mathcal{E}$ be a sequence converging in variation to a distribution P . We need to show that P also satisfies the linear constraints and has uniform marginals. Convergence in variation implies that for every test function f

$$\int_{[0,1]^d} f(u) p_n(u) du \rightarrow \int_{[0,1]^d} f(u) p(u) du. \quad (13)$$

Taking $f = \phi_m$ proves that the limit distribution P satisfies the linear constraints. We now need to prove that it has uniform marginals.

Let us consider $u_{-i} = (\dots, u_{i-1}, u_{i+1}, \dots)$ the vector u without its i -th coordinate, and let us choose a test function f that only depends on u_i . By Fubini’s theorem we have

$$\begin{aligned} \int_{[0,1]^d} f(u) p_n(u) du &= \int f(u_i) \underbrace{\left(\int p_n(u) du_{-i} \right)}_{=1} du_i \\ &= \int_{[0,1]} f(u_i) du_i \\ &= \int_{[0,1]^d} f(u_i) p(u) du \end{aligned}$$

where the last equality is due to the Equation (13). Putting the last two equality together, we get

$$\int f(u_i) \left(1 - \int p(u) du_{-i} \right) du_i = 0$$

for every univariate test function f , which implies that every marginal of P is uniform. \square

\mathcal{E} being convex, non-empty, and variation closed, U admits a unique I -projection on \mathcal{E} (see Theorem 2.1

in Csiszár (1975)), or equivalently, (MIND) admits a unique solution.

Functional form of the pdf: Let us denote \mathcal{F} the space of distributions supported on $[0, 1]^d$ that satisfy the linear constraint $E_P[\phi_m(\mathbf{u})] = \alpha_m$. Clearly, $\mathcal{E} \subset \mathcal{F}$ as the only difference between the two sets is that \mathcal{E} only contains elements of \mathcal{F} with uniform marginals (i.e. copula distributions). The I -projection P_{AM} of U on \mathcal{F} , which exists because \mathcal{F} is convex, non-empty, and variation closed, is the minimizer of the problem (A-MIND).

By Theorem 2.3 in Csiszár (1975), P_M is the I -projection of P_{AM} on \mathcal{E} .

Applying Theorem 3.1 (Case A) in Csiszár (1975), we get that P_{AM} has density with respect to U , which is also its pdf, of the form $p_{AM} = e^{\theta^T \phi_m(\mathbf{u})}$. Moreover, any distribution in \mathcal{F} with density with respect to U of this form is the I -projection of U on \mathcal{F} .

Applying Theorem 3.1 (Case B) in Csiszár (1975), we get that P_M has density with respect to P_{AM} of the form $\prod_{i=1}^d f_i(u_i)$, where f_i are non-negative and log-integrable. Moreover, any distribution in \mathcal{E} with density with respect to P_{AM} of this form is the I -projection of P_{AM} on \mathcal{E} .

Hence, P_M has density with respect to U , which is also its pdf,

$$p_M = e^{\theta^T \phi_m(\mathbf{u})} \prod_{i=1}^d f_i(u_i),$$

and any distribution on $[0, 1]^d$ whose pdf of this form is the minimizer of (MIND).

Pythagoras' Identity: Theorem 3.1 in Csiszár (1975) guarantees that identity (3.1) in Csiszár (1975) holds and

$$-h(P) = -h(P_{AM}) + KL(P||P_{AM}) \quad (14)$$

for any $P \in \mathcal{F}$ and

$$KL(Q||P_{AM}) = KL(P_M||P_{AM}) + KL(Q||P_M) \quad (15)$$

for any $Q \in \mathcal{E} \subset \mathcal{F}$.

As $Q \in \mathcal{F}$,

$$KL(Q||P_{AM}) = h(P_{AM}) - h(Q).$$

Thus,

$$h(P_{AM}) - h(Q) = KL(P_M||P_{AM}) + KL(Q||P_M)$$

and

$$-h(Q) = -h(P_{AM}) + KL(P_M||P_{AM}) + KL(Q||P_M).$$

As $P_M \in \mathcal{F}$,

$$-h(P_M) = -h(P_{AM}) + KL(P_M||P_{AM}),$$

and we get

$$-h(Q) = -h(P_M) + KL(Q||P_M).$$

B.3 Proof of Theorem 3.2

To prove Theorem 3.1 we had to prove Theorem 3.2. See Section B.2.

B.4 Proof of Theorem 3.3

Let

$$g(\mathbf{u}; \phi_m, \alpha_m) := \prod_{i=1}^q p_M(w, \mathbf{v}_i; \eta_i, \beta_i).$$

We want to prove that

$$g(\mathbf{u}; \phi_m, \alpha_m) = p_M(\mathbf{u}; \phi_m, \alpha_m).$$

We know from Theorem 3.1 that $p_M(\mathbf{u}; \phi_m, \alpha_m)$ is the only copula density of the form

$$e^{\theta^T \phi_m(\mathbf{u})} = e^{\sum_{i=1}^q \theta_i^T \eta_i(w, \mathbf{v}_i)} = \prod_{i=1}^q e^{\theta_i^T \eta_i(w, \mathbf{v}_i)}$$

that satisfies the constraints $E_P[\psi_m(\mathbf{u})] = \beta_m$.

First we note that g is a copula entropy. Indeed, integrating g with respect to every variable but a coordinate of \mathbf{v}_i is always 1 by virtue of the fact that $p_M(w, \mathbf{v}_i; \eta_i, \beta_i)$ are copula densities. To see why, note that we may first integrate with respect to \mathbf{v}_j for all $j \neq i$, and then with respect to w and all other coordinates of \mathbf{v}_i . Additionally, if we integrate with respect to all variables but w , we get

$$\begin{aligned} & \int g(\mathbf{u}; \phi_m, \alpha_m) d\mathbf{v}_1 \dots d\mathbf{v}_q \\ &= \int p_M(w, \mathbf{v}_q; \eta_q, \beta_q) \times \dots \times \\ & \quad \left(\underbrace{\int p_M(w, \mathbf{v}_1; \eta_1, \beta_1) d\mathbf{v}_1}_{=1} \right) d\mathbf{v}_2 \dots d\mathbf{v}_q \\ &= 1. \end{aligned}$$

Second, g clearly has the form $\prod_{i=1}^q e^{\theta_i^T \eta_i(w, \mathbf{v}_i)}$.

Finally, g satisfies the constraints $E_P[\psi_m(\mathbf{u})] = \beta_m$

as

$$\begin{aligned}
 & \int \eta_i(w, \mathbf{v}_i) g(\mathbf{u}; \phi_m, \alpha_m) d\mathbf{u} \\
 &= \int \eta_i(w, \mathbf{v}_i) p_M(w, \mathbf{v}_i; \eta_i, \beta_i) \\
 & \quad \times \left(\underbrace{\int \prod_{j \neq i} p_M(w, \mathbf{v}_j; \eta_j, \beta_j) d\mathbf{v}_j}_{=1} \right) d\omega d\mathbf{v}_i \\
 &= \int \eta_i(w, \mathbf{v}_i) p_M(w, \mathbf{v}_i; \eta_i, \beta_i) d\omega d\mathbf{v}_i \\
 &= \beta_i,
 \end{aligned}$$

where we've used the fact that each $p_M(w, \mathbf{v}_j; \eta_j, \beta_j)$ has uniform marginals, and satisfies the constraint

$$E_P[\eta_j(w, \mathbf{v}_j)] = \beta_j.$$

B.5 Proof of Theorem 3.4

The essence of the proof is in the transitivity property of I -projections. Indeed, if $\mathcal{P} \subset \mathcal{Q}$ are linear sets of probability distributions supported on $[0, 1]^d$, Q the I -projection of the standard uniform U on \mathcal{Q} , and P the I -projection of U on \mathcal{P} , then P is also the I -projection of Q on \mathcal{P} (Theorem 2.3 Csiszár (1975)).

In the case of Theorem 3.4, \mathcal{Q} is the set of probability distributions satisfying all constraints, and \mathcal{Q} is the set of probability distributions satisfying all but the between-blocks constraints. If we denote, $q(\mathbf{u}; \psi_m, \beta_m) := \prod_{i=1}^q p_{AM}(\mathbf{v}_i; \eta_i, \beta_i)$, a direct application of Theorem 4.1 shows that q is the density of the I -projection Q of U on \mathcal{Q} . The maximizer of the full (A-MIND) problem is therefore the I -projection of Q on \mathcal{P} , and we know from Theorem 3.1 in Csiszár (1975) that it has Radon-Nikodym derivative with respect to Q of the form $\frac{dP}{dQ} = e^{\theta^T \phi_m^k(\mathbf{u})}$. Putting everything together, we get that the maximizer P of the full (A-MIND) problem has pdf of the form

$$p_{AM}(\mathbf{u}; \phi_m, \beta_m) = e^{\theta^T \phi_m^k(\mathbf{u})} \prod_{i=1}^q p_{AM}(\mathbf{v}_i; \eta_i, \beta_i).$$

The unicity of this representation is a direct consequence of Theorem 4.1. Taking the negative log of this expression and then the expectation, we get

$$\begin{aligned}
 h_{AM}(\mathbf{u}; \phi_m^k, \beta_m) &= -\theta^T \alpha_m^k \\
 &\quad - \sum_{i=1}^q E_{p_{AM}(\mathbf{v}_i; \phi_m^k, \beta_m)} [\log p_{AM}(\mathbf{v}_i; \eta_i, \beta_i)].
 \end{aligned}$$

The final result stems from the identity

$$E_P(-\log q) = E_P(-\log p) - \text{KL}[p||q].$$

The fact that

$$-\theta^T \alpha_m^k \leq \sum_{i=1}^q \text{KL}[p_{AM}(\mathbf{v}_i; \phi_m^k, \beta_m) || p_{AM}(\mathbf{v}_i; \eta_i, \beta_i)]$$

is a direct consequence of $\mathcal{P} \subset \mathcal{Q}$, which implies that the entropy of P cannot be greater than that of Q , and the fact that $\sum_{i=1}^q h_{AM}(\mathbf{v}_i; \eta_i, \beta_i)$ is the entropy of Q . The two entropies are the same if and only if $P = Q$ or, equivalently, $\theta = 0$. When $\forall \mathbf{u}$, $\gamma(\mathbf{u})$ is constant and equal to 1, $P = Q$.

B.6 Proof of Lemma 4.1

The Hessian of the objective, namely $\int_{[0,1]^d} (\phi_m(\mathbf{u}) \phi_m(\mathbf{u})^T) e^{\theta^T \phi_m(\mathbf{u})} d\mathbf{u}$, is clearly strictly positive-definite as coordinates of 1 are not linearly related.

B.7 Proof of Lemma 4.2

The only critical point of the objective of (CVX-MIND) satisfies

$$\alpha_m = \int_{[0,1]^d} \phi_m(\mathbf{u}) e^{\theta^{*T} \phi_m(\mathbf{u})} d\mathbf{u}.$$

Given that the first coordinates of ϕ_m and α_m are both 1, $\int_{[0,1]^d} e^{\theta^{*T} \phi_m(\mathbf{u})} d\mathbf{u} = 1$. By Theorem 3.2, the distribution with pdf $e^{\theta^{*T} \phi_m(\mathbf{u})}$ maximizes (A-MIND).

B.8 Proof of Theorem 4.1

Theorem 4.1-A is a consequence of the uniqueness of the solution to the Hausdorff moment problem.

Indeed, the uniform distribution on $[0, 1]$ is uniquely characterized by the sequence of moments $\forall j, E(u^j) = 1/(1+j)$ (Shohat and Tamarkin (1943)). Thus, we may replace the uniform marginal constraints in (MIND) with the constraints $\forall i, j, E(u_i^j) = 1/(1+j)$. The only difference between (A-MIND) and (MIND) is that the former has k moment constraints whereas the latter has all moment constraints. It follows that

$$\forall m > 0, \quad h_{AM}(\mathbf{u}; \phi_m^k, \beta_m) \xrightarrow{k \rightarrow \infty} h_M(\mathbf{u}; \psi_m, \beta_m).$$

Additionally, a direct application of the basic consistency theorem for extremum estimators (see Newey and McFadden (1994)) to (CVX-MIND) shows that for any consistent estimator $\hat{\beta}_{m,n}$ of β_m ,

$$h_{AM}(\mathbf{u}; \phi_m^k, \hat{\beta}_{m,n}) \xrightarrow{n \rightarrow \infty} h_{AM}(\mathbf{u}; \phi_m^k, \beta_m).$$

As for Theorem 4.1-B, we simply need to prove that:

$$\forall k > 0, \quad h_{AM}(\mathbf{u}; \phi_m^k, \beta_m) \xrightarrow{m \rightarrow \infty} h(\mathbf{u}_z).$$

We recall that $h(\mathbf{u}_z) = -KL(P_{\mathbf{u}_z}||U)$. Using the NWJ characterization of the KL divergence (Nguyen et al. (2010)), we get:

$$h(\mathbf{u}_z) = \inf_{T \in \mathcal{T}} -E_{P_{\mathbf{u}_z}}[T(\mathbf{u})] + \int_{[0,1]^d} e^{T(\mathbf{u})} - 1 d\mathbf{u}, \quad (16)$$

where \mathcal{T} is the space of continuous functions on $[0, 1]^d$. Using the fact that $(\phi_m)_m$ is dense in \mathcal{T} (property (P3)), we may rewrite Equation (16) as

$$\begin{aligned} h(\mathbf{u}_z) = \lim_{m \rightarrow \infty} \min_{\boldsymbol{\theta}_m} & -1 - \underbrace{\boldsymbol{\theta}_m^T E_{P_{\mathbf{u}_z}}[\phi_m(\mathbf{u})]}_{:= \alpha_m} \\ & + \int_{[0,1]^d} e^{\boldsymbol{\theta}_m^T \phi_m(\mathbf{u})} d\mathbf{u}, \end{aligned} \quad (17)$$

where $\boldsymbol{\theta}_m$ satisfies $\boldsymbol{\theta}_m^T \phi_m(\mathbf{u}) = T(\mathbf{u}) - 1$. Note that the inner optimization problem has the same minimizer as (CVX-MIND), and the minimum is $h_{\text{AM}}(\mathbf{u}; \phi_m^k, \boldsymbol{\beta}_m)$, where we have used the fact that the first coordinate of ϕ_m is 1.

Putting everything together, we get

$$\forall k, h(\mathbf{u}_z) = \lim_{m \rightarrow \infty} h_{\text{AM}}(\mathbf{u}; \phi_m^k, \boldsymbol{\beta}_m).$$